# ABSTRACT HIDDEN MARKOV MODELS:
# A MONADIC ACCOUNT OF QUANTITATIVE INFORMATION FLOW

ANNABELLE MCIVER $^a$, CARROLL MORGAN $^b$, AND TAHIRY RABEHAJA $^a$

$^a$ Dept. Computing, Macquarie University., Sydney, Australia
   *e-mail address*: {annabelle.mciver,tahiry.rabehaja}@mq.edu.au

$^b$ School of Comp. Sci. and Eng., Univ. New South Wales, and Data61., Sydney, Australia
   *e-mail address*: carroll.morgan@unsw.edu.au

ABSTRACT. Hidden Markov Models, *HMM*'s, are mathematical models of Markov processes with state that is hidden, but from which information can leak. They are typically represented as 3-way joint-probability distributions.

We use *HMM*'s as denotations of probabilistic hidden-state sequential programs: for that, we recast them as "abstract" *HMM*'s, computations in the Giry monad $\mathbb{D}$, and we equip them with a partial order of increasing security. However to encode the monadic type *with hiding* over some state $\mathcal{X}$ we use $\mathbb{D}\mathcal{X} \to \mathbb{D}^2\mathcal{X}$ rather than the conventional $\mathcal{X} \to \mathbb{D}\mathcal{X}$ that suffices for Markov models whose state is not hidden. We illustrate the $\mathbb{D}\mathcal{X} \to \mathbb{D}^2\mathcal{X}$ construction with a small Haskell prototype.

We then present *uncertainty measures* as a generalisation of the extant diversity of probabilistic entropies, with characteristic analytic properties for them, and show how the new entropies interact with the order of increasing security. Furthermore, we give a "backwards" uncertainty-*transformer* semantics for *HMM*'s that is dual to the "forwards" abstract *HMM*'s — it is an analogue of the duality between forwards, relational semantics and backwards, predicate-transformer semantics for imperative programs with demonic choice.

Finally, we argue that, from this new denotational-semantic viewpoint, one can see that the Dalenius desideratum for statistical databases is actually an issue in compositionality. We propose a means for taking it into account.

## 1. INTRODUCTION

1.1. **Setting and overview.** We can represent probabilistic sequential programs with hidden state as Hidden Markov Models, i.e. *HMM*'s [1] formulated as probabilistic mechanisms that take prior, input probability distributions and give posterior distributions over (leaked) observations and final state. Here, however, we recast *HMM*'s as computations over the Giry monad, making them more suitable for denotational semantics. Indeed the monadic view of simple Markov processes in particular is well established [1, 2], using $\mathcal{X} \to \mathbb{D}\mathcal{X}$ where

---

[1]We use apostrophe uniformly between acronyms and suffixes, even when they are not possessive.

type-constructor $\mathbb{D}$ makes distributions on its base type $\mathcal{X}$; the Kleisli extension is then of type $\mathbb{D}\mathcal{X} \to \mathbb{D}\mathcal{X}$, representing the action of multiplying an initial-state-distribution vector by a Markov matrix. But that simplicity cannot account for hidden state and information flow.

We treat hidden state by beginning with $\mathbb{D}\mathcal{X}$ (not $\mathcal{X}$): the computation type we obtain is then "one level up", of type $\mathbb{D}\mathcal{X} \to \mathbb{D}^2\mathcal{X}$, the Kleisli extension is $\mathbb{D}^2\mathcal{X} \to \mathbb{D}^2\mathcal{X}$; and we call the double-distribution type $\mathbb{D}^2\mathcal{X}$ *hyper-distributions*, or "hypers" for short.

Although the Giry monad is formulated in terms of general measures [2], we will need only discrete distributions for matrix-based *HMM*'s. Nevertheless, we give our constructions and results in more general terms, anticipating e.g. infinite sequences of *HMM*'s, nondeterminism, and iterations for which proper measures will be necessary [3].

In earlier work, we have used the hypers $\mathbb{D}^2\mathcal{X}$, equipped with a partial order of increasing security, to establish compositionality results [4], to explore the effect of including demonic nondeterminism [5] and to give an abstract treatment of probabilistic channels [6, 7]. A second earlier theme has been the generalisation of entropies (such as Shannon) to a more abstract setting where only their essential properties are preserved [4, 5, 7, 8]. Here we use monads to bring all those separate strands together and to go further.

One major further step is to show that there is a dual, backwards view for abstract *HMM*'s, based on "uncertainty" *transformers* that transform post- uncertainty measures into pre- uncertainty measures where, in turn, *uncertainty measures* generalise probabilistic entropies.

We and others have argued that specific entropies (e.g. Shannon) have limitations in security work generally [4, 9]. Therefore we focus here on their essential properties: continuity and concavity. That view is supported by powerful theorems that such a generalisation supports, and a methodological criterion that uncertainty measures capture contexts in a way that individual styles of entropy cannot.

A second further step is to extend our recent treatment [7] of the *Dalenius Desideratum*, the "collateral" leakage of information due to unknown correlations with third-party data, from merely channels (a "read only" scenario [10, 11], such as access to a statistical database) to programs that might alter the database (thus "read/write" as well). The Dalenius perspective here is the fact that care must be taken wrt. compositionality in a context containing extra variables even when a program fragment does not explicitly refer to them[4].

*To remain accessible to a broader security community*, we do not begin from Giry: rather we first work in elementary terms. In §7 the monadic structures will be seen to have informed our earlier definitions and theorems.

## 1.2. **Principal contributions and aims: summary.**
Our principal contributions are these, in which the **new constructions and results** are given in bold:

- We note that (finite) classical *HMM*'s are a model for straight-line sequential probabilistic programs with hidden state
- We formulate *abstract HMM's* over a state as a monadic model for *HMM*'s over that same state, and give their characteristic properties.
- We formulate *uncertainty measures* as a generalisation of diverse entropies (top centre), and give their characteristic properties. [2]
- We note that uncertainty measures have a complete representation based on real-valued functions of state and adversarial strategy.

---

[2]They were studied, but less extensively, as "disorders", in [5].

- We give a dual, uncertainty-transformer semantics of *HMM*'s and prove the duality.
- We show how all of the above is an instance of the general Giry monad as a computation, of which (finite) *HMM*'s use a discrete portion.
- We explain how the "Dalenius effect" is manifested as a compositional issue in this framework, and how it can be treated.

In other sections we review abstract channels (§2.2), hyper-distributions (§2.3) and the security order (§6) on hypers.

We believe that Thm. 9.8, in particular its assumptions and proof, is a significant new result.

Our principal <u>aims</u> are these:

- (More abstract) To construct forward- and dual backward semantic spaces for probabilistic sequential computations over hidden state, using monadic computations and partial (refinement) orders in this new context, and we formulate and prove the general properties that make them suitable for embedding finite (for the moment) *HMM*'s.
- (More concrete) To provide the basis for a source-level reasoning method, analogous to Hoare logic or weakest preconditions, for quantitative non-interference in sequential programs. For this, the dual, transformer semantics for *HMM*'s seems to be a necessary first step, together with a link between the social aspects of security and the mathematical behaviour of a program (§11).

The conclusion §14 discusses the <u>benefits</u> of doing this.

### 1.3. General notations — see also §A.
Application of function $f$ to argument $x$ is written $f.x$ to reduce parentheses. It associates to the left.

Although a matrix $M$ with rows, columns indexed by $R, C$ is a function $R \times C \to \mathbb{R}$, we avoid constant reference to the reals $\mathbb{R}$ by writing just $R \twoheadrightarrow C$ for that type; similarly we write the type of a vector over $X$ as $\overrightarrow{X}$. We write $M_{r,c}$ for the element of matrix $M$ indexed by row $r$ and column $c$; then the $r$-th row of $M$ is $M_{r,-}$; and the $c$-th column is $M_{-,c}$, of types $\overrightarrow{Y}, \overrightarrow{X}$ resp. For row- or column vector $v \colon \overrightarrow{I}$ we write $v_i$ for its $i$-th element. Thus e.g. we have $(M_{-,c})_r = M_{r,c}$.

When multiplying vectors and matrices we assume without comment that the vector has been oriented properly, i.e. as a row or column as required. Thus $v$ acts as a row in $v \cdot M$ but as a column in $M \cdot v$. Thus for $v \colon \overrightarrow{X}$ and $M \colon X \twoheadrightarrow Y$ the matrix product $v \cdot M$ is in $\overrightarrow{Y}$, where here we are using dot ($\cdot$) for matrix multiplication. Multiplication of scalars will usually be juxtaposition, but occasionally $\times$ when we are avoiding ambiguity.

We write for example $x \colon X$, i.e. with a colon, when we are introducing a fresh variable $x$ into the discussion at that point; with $x \in X$ we are instead stating a property of some $x$ and $X$ that have been already introduced at some earlier point. [3] That means in the former case that one need not search backwards to see what $x$ is being referred to (and in the latter case, one might).

Other specific notations are explained at first use, and (as noted above) a full glossary in occurrence order is given in §A.

---

[3] For example we could write "Because we have already established that $s \in \mathbb{P}X$, we know that for any $x \colon s$ we have $x \in X$." Both $s, X$ are defined in the surrounding text, but $x$ here is a local (i.e. bound) variable just used temporarily.

## 2. Abstract channels and hyper-distributions

We now review abstract channels as a conceptual stepping-stone to hyper-distributions — recall they are "hypers" for short. (Channels are the special case of *HMM*'s where the state is not updated.)

### 2.1. Channels and distributions as matrices and vectors.
A *channel* is a (stochastic) matrix of non-negative reals with 1-summing rows; we use upper-case Roman letters like $C$ for them. The rows are labelled with elements from some set $\mathcal{X}$; and the columns from some set $\mathcal{Y}$. Thus a channel typically has type $\mathcal{X} \rightarrow \mathcal{Y}$; here, both $\mathcal{X}$ and $\mathcal{Y}$ will be finite.

A distribution in $\mathbb{D}\mathcal{X}$ can be presented as a 1-summing vector in $\overrightarrow{\mathcal{X}}$, usually lower-case Greek: generally $\delta$ for "distribution", but especially $\pi$ for prior and sometimes $\rho$ for posterior.

**Definition 2.1** (Weight). Let $M$ or $v$ be a matrix or vector resp. Then $\Sigma M$ or $\Sigma v$ is its *weight*, the sum $\Sigma_{x,y} M_{x,y}$ or $\Sigma_x v_x$ taken over all its indices.

Thus e.g. we have $\Sigma M_{x,-} = \Sigma_y M_{x,y}$ and that $M$ is stochastic (i.e. represents a channel) just when $\Sigma M_{x,-}$ is 1 for all $x$.

Each row $C_{x,-}$ of a channel $C$ is a conditional probability distribution over $\mathcal{Y}$ given that particular $x{:}\mathcal{X}$. That is, the $y$-th element $C_{x,y}$ of $C_{x,-}$ is the probability that $C$ takes input $x$ to output $y$.

### 2.2. Informal channel semantics: *abstract* channels.
A (1-summing) prior $\pi$ and (stochastic) channel $C$ together determine a joint distribution as follows.

**Definition 2.2** (Channel applied to prior). Given a prior $\pi{:}\overrightarrow{\mathcal{X}}$ and channel $C{:}\mathcal{X} \rightarrow \mathcal{Y}$ we write $\pi \triangleright C$ for the joint-distribution matrix of type $\mathcal{X} \rightarrow \mathcal{Y}$ resulting from *applying* the channel to the prior, defined $(\pi \triangleright C)_{x,y} := \pi_x C_{x,y}$. (Here juxtaposition is ordinary multiplication of reals.)

Note that matrix $\pi \triangleright C$ is not stochastic: rather because $C$ itself is stochastic we have $\Sigma(\Sigma(\pi \triangleright C)_{x,-}) = \Sigma \pi_x = 1$.

A non-zero vector is normalised as follows.

**Definition 2.3** (Normalisation). Let $\delta{:}\overrightarrow{\mathcal{X}}$ be such that $0 \neq \Sigma \delta$. Then the *normalisation* $\lfloor \delta \rfloor$ of $\delta$ is given by $\lfloor \delta \rfloor_x := \delta_x / \Sigma \delta$ for each $x{:}\mathcal{X}$.

Now for some $\pi{:}\overrightarrow{\mathcal{X}}$ and channel $C{:}\mathcal{X} \rightarrow \mathcal{Y}$ define joint distribution $J{:}\mathcal{X} \rightarrow \mathcal{Y}$ by $J = \pi \triangleright C$. The (marginal) probability of each output $y{:}\mathcal{Y}$ is $\Sigma J_{-,y}$ and, associated with each $J$, there is a posterior distribution $\lfloor J_{-,y} \rfloor$ on $\mathcal{X}$.

Abstracting from the $y$-values, but retaining the link between the marginal probabilities and the posterior distributions, gives an informal description of our intended "abstract channel" semantics [6]. We make this precise in §2.4.

2.3. **Hypers abstract from joint distributions.** The joint-distribution matrix $J=\pi\triangleright C$ contains "too much" information if we do not need the actual value of $y$ that led to a particular posterior. Abstracting from those output values leads us to a representation of the possible posteriors on their own, retaining however the probability with which they occur (in fact the marginal probability of the value $y$ that produced each one). The advantage we take from that is that *HMM*'s acquire a monadic structure, acting as Kleisli maps, and furthermore can express other probabilistic notions in a way more suited to calculation: for example, conditional entropies become expected values (of the entropies) over the distribution of posteriors.

More intuitive reasons for the abstraction include that it is appropriate in security to consider the information leakage of a channel $C$ wrt. a prior $\pi$ to concern only what an adversary can discover about $\pi$, and not the actual observations that led to that discovery: whether a spy's vocabulary is "da/nyet" or "yes/no", or indeed whether "yes" means "it's zero" or "it's one", does not affect the information-theoretic threat that spy represents, provided of course that the spy and her controller have agreed on the vocabulary beforehand.

We can abstract from the observations in $\pi\triangleright C$ as follows. If column $y$ of $J = \pi\triangleright C$ is all zero, then that $y$ will never occur (for any prior); thus we can omit that column.

And if two columns $y_{1,2}$ of $J$ are proportional to each other, i.e. are *similar* (as for triangles), then we can add them together, since for a given prior the same posterior will be inferred for $y_1$ as for $y_2$ and the overall probability of inferring that posterior will be the sum of the marginal probabilities for $y_{1,2}$.[4]

Finally, a 1-1 renaming of the $y$-values has no effect on the posteriors and their respective probabilities; so we can remove those names as long as we retain the distinction between separate (non-zero, non-similar) columns.

Abstracting from all that arguably inessential information (about $y$) leaves only a distribution of posteriors on $\mathcal{X}$ and, for us, this is the semantic view. Writing in general $\mathbb{D}\mathcal{X}$ for 1-summing functions of type $\mathcal{X}\to\mathbb{R}^{\geq}$, a discrete distribution over $\mathcal{X}$ has type $\mathbb{D}\mathcal{X}$ and so a discrete distribution of such distributions has type $\mathbb{D}(\mathbb{D}\mathcal{X})$ that is $\mathbb{D}^2\mathcal{X}$. Those latter are our hypers, and they are our abstraction of joint distributions $\mathcal{X}\twoheadrightarrow\mathcal{Y}$.

The values of type $\mathbb{D}\mathcal{X}$ are called the *inners* of a hyper, and the *outer* distribution of a hyper is its distribution over those inners: that is, a hyper on $\mathcal{X}$ is a (single) outer distribution over (possibly many) inners, and each inner is a (single) distribution over $\mathcal{X}$ itself.

As an example, recall the famous puzzle of *Bertrand's Boxes*. Three identical boxes contain two balls each: one has two white balls; one has two black balls; and the remaining box has one of each. It is not known which box is which; and one of them is chosen randomly. A ball is drawn at random from it, and it is white. What is the probability that the other ball in that box is also white? We reason as follows.

The state space is $\mathcal{X}=\{0,1,2\}$, referring to the number of white balls in each box. The prior distribution in $\mathbb{D}\mathcal{X}$ is uniform, which we can write $(1/3, 1/3, 1/3)$. The *HMM* is a channel that takes input $x$ to the distribution $(white \mapsto x/2, black \mapsto 1-x/2)$. The joint distribution $p$ say, of type $\mathbb{D}(\mathcal{X}\times\{white, black\})$, would be such that $p(1, white) = 1/3\times1/2 = 1/6$, the probability that Box 1 was chosen and that the ball taken from it was white. The overall probability that a white ball is taken (from whichever box) is the *white*-marginal $0+1/6+1/3 = 1/2$ (which is obvious from symmetry anyway), and the posterior distribution

---

[4]For brevity we write $y_{1,2}$ rather than $y_1, y_2$.

on $\mathcal{X}$ is in that case $(0, 1/3, 2/3)$ — which nicely solves the puzzle. The posterior probability that $x{=}2$ is $2/3$ given that a white is taken and, by the way, a white is taken with overall probability $1/2$ (the marginal). And for that reasoning of course the value of the observation, the colour of the drawn ball, is used.

But now suppose instead you wanted to know only the decrease in Shannon entropy resulting from that experiment. Beforehand the entropy is $H(1/3, 1/3, 1/3) = \log_2 3 = 1.58$ (approximately). Afterwards, it will be the *conditional* Shannon entropy of the distribution of posteriors, calculated by taking the expected value of $H()$ over the distribution of posteriors: and that is approximately

$$1/2{\times}H(0, 1/3, 2/3) + 1/2{\times}H(2/3, 1/3, 0) \quad = \quad 1/2{\times}0.92 + 1/2{\times}0.92 \quad = \quad 0.92 \quad , \qquad (2.1)$$

so that $1.58{-}0.92 = 2/3$ (exactly) of a bit has been leaked. And we did not need colours for that: the calculation is done entirely with the hyper-distribution, that is with the distribution

$$\text{the } \textit{outer} \left\{ \begin{array}{ccc} 1/2 & \mapsto & (0, 1/3, 2/3) \\ 1/2 & \mapsto & (2/3, 1/3, 0) \end{array} \right\} \text{the } \textit{inners}$$

*of distributions*, i.e. of posteriors: a hyper-distribution. The $1/2$'s are the marginals, and the $(\cdots)$ are the posteriors associated with each. We call the channel-output marginal the *outer*, and the posteriors the *inners*.

Seeing this example as a security leak, we might imagine that the adversary is trying to guess the colour of the other ball in the box: in that case she would look at the colour she took and then guess that same colour. To describe that we use a different entropy $V_1$, called *Bayes Vulnerability*, which is the probability the secret can be guessed in one try by an optimal adversary. Obviously she will guess the $x$-value with the largest probability in the posterior (the inner), and her conditional probability of guessing correctly is

$$1/2{\times}V_1(0, 1/3, 2/3) + 1/2{\times}V_1(2/3, 1/3, 0) \quad = \quad 1/2{\times}2/3 + 1/2{\times}2/3 \quad = \quad 2/3 \quad .$$

That's no surprise — but what is worth noting is that we used the *same* hyper-distribution for the $V_1$ calculation just above as for the $H$ calculation at (2.1). That is the utility of the abstraction: that the hyper contains enough information to handle many entropies one might use to measure leakage.

2.4. **The semantic function from joints to hypers.** In this section we define precisely the denotation $[\![J]\!]$ in $\mathbb{D}^2\mathcal{X}$ of a joint-distribution matrix $J\colon \mathcal{X}{\rightarrow}\mathcal{Y}$.

**Definition 2.4** (Sub-distribution, sup-hyper)**.** A discrete sub-distribution over a set $\mathcal{X}$ is a function of type $\mathcal{X}{\rightarrow}[0,1]$ that sums to *no more than 1*; we write that type as $\underline{\mathbb{D}}\mathcal{X}$. (Recall that a *proper* distribution in $\mathbb{D}\mathcal{X}$ sums to exactly 1, and thus $\mathbb{D}\mathcal{X} \subseteq \underline{\mathbb{D}}\mathcal{X}$.)

Similarly a discrete sub-*hyper* over a set $\mathcal{X}$ is a sub-distribution over the (proper, inner) distributions $\mathbb{D}\mathcal{X}$, thus of type $\underline{\mathbb{D}}(\mathbb{D}\mathcal{X})$; only the outer of a sub-hyper can sum to less than 1. We write that type as $\underline{\mathbb{D}}^2\mathcal{X}$. (Note that the inners of a sub-hyper are proper distributions.)

**Definition 2.5** (One- and two-point distributions)**.** For $z, z'\colon \mathcal{Z}$ in general we write $[z]$ for the *point distribution* on $z$, viz. assigning probability 1 to $z$ and 0 to all other elements of $\mathcal{Z}$.[5] We write $z_p{\oplus}z'$ for the two-point distribution that assigns $p$ to $z$ and $1{-}p$ to $z'$ and 0 to everything else in $\mathcal{Z}$.

Thus $z_{\,1}{\oplus}\, z' = [z]$ and $z_{\,0}{\oplus}\, z' = [z']$.

---

[5]Function $[\cdot]$ is the unit $\eta$ of the $\mathbb{D}$-monad: see §7.

**Definition 2.6** (Point sub-hyper). For sub-distribution $\delta\colon\mathbb{D}\mathcal{X}$ the *point sub-hyper* $[\delta]$ in $\mathbb{D}^2\mathcal{X}$ has weight $\Sigma\delta$ concentrated on the single (inner) distribution $\lfloor\delta\rfloor$, provided of course that $\Sigma\delta\neq 0$. If $\Sigma\delta=0$ then $[\delta]$ is the (unique) weight-zero sub-hyper.

That is, the argument $\delta$ is normalised to make the inner, and its weight becomes the (one-point) sub-outer on that inner.

Note that when $\Sigma\delta = 1$ we have $[\delta] = [\delta]$.

We now define the semantic function itself:

**Definition 2.7** (Joint-distribution denotes hyper). Let $J\colon\mathcal{X}\rightarrowtail\mathcal{Y}$ satisfy $1=\Sigma J$ so that it describes a discrete (proper) joint distribution in $\mathbb{D}(\mathcal{X}\times\mathcal{Y})$. Then its abstraction $[\![J]\!]$ to a hyper in $\mathbb{D}^2\mathcal{X}$ is given by

$$[\![J]\!] \quad = \quad \sum_{y\in\mathcal{Y}}[J_{-,y}] \quad ,$$

with summation $\Sigma_{y\in\mathcal{Y}}$ therefore being an addition of sub-hypers, i.e. sub-distributions on $\mathbb{D}\mathcal{X}$. Each column $J_{-,y}$ is regarded as a sub-distribution in $\underline{\mathbb{D}\mathcal{X}}$, and then $\underline{[-]}$ converts it to a sub-point hyper.

Note that in Def. 2.7 any all-zero columns in $J$ are automatically ignored, since they become zero-weight sub-hypers in the sum and drop out automatically. If however all columns of $J$ are zero, then its denotation $[\![J]\!]$ becomes automatically the weight-zero sub-point hyper.

2.5. **Abstract channels — review.** In earlier work [6] we described an "abstract channel" as a function from prior distributions to hypers. We restate that here in our current denotational style:

**Definition 2.8** (Denotation of channel). Let $C\colon\mathcal{X}\rightarrowtail\mathcal{Y}$ be a channel matrix. Its denotation, of type $\mathbb{D}\mathcal{X}\to\mathbb{D}^2\mathcal{X}$, is called an *abstract channel* and is defined for $\pi\colon\mathbb{D}\mathcal{X}$ by

$$[\![C]\!].\pi:=[\![\pi\triangleright C]\!] \ ,$$

where the $[\![\cdot]\!]$ on the left is the denotational function for channels, being defined here; and on the right it is the denotational function on joint distributions that we have already from Def. 2.7. (We use $[\![-]\!]$ uniformly for denotation functions, relying on context instead of e.g. using subscripts like $[\![-]\!]_{\text{chan}}$ on the left and $[\![-]\!]_{\text{joint}}$ on the right.)
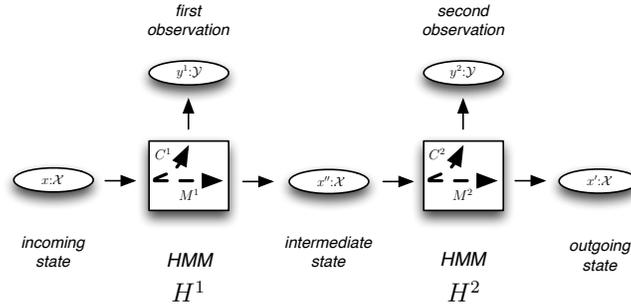
In fact the prior $\pi$ can be recovered from $\pi\triangleright C$, as this definition shows:

**Definition 2.9** (Support of a distribution). Given discrete distribution $\delta\colon\mathbb{D}\mathcal{Z}$, we write $\lceil\delta\rceil$ for the *support* of $\delta$, the set of elements $z\colon\mathcal{Z}$ for which $\delta.z$, the probability assigned by $\delta$ to $z$, is not zero. Obviously $\delta\in\mathbb{D}\mathcal{Z}$ implies $\lceil\delta\rceil\subseteq\mathcal{Z}$; if in fact $\lceil\delta\rceil=\mathcal{Z}$ then we say that $\delta$ is *full support*.

**Definition 2.10** (Average of a hyper). For hyper $\Delta\colon\mathbb{D}^2\mathcal{X}$ define its average $\mathsf{avg}.\Delta$ in $\mathbb{D}\mathcal{X}$ by

$$\mathsf{avg}.\Delta.x \quad := \quad \sum_{\delta\colon\lceil\Delta\rceil}(\Delta.\delta)(\delta.x) \qquad \text{for all } x\colon\mathcal{X}, \quad {}^6$$

where we use upper-case Greek for hypers.

Each step $H^{1,2}$ takes an input- to an output state in $\mathcal{X}$; the observations $y^{1,2}$: $\mathcal{Y}$ are accumulated. In each step $H^{1,2}$ the output state is determined by a markov $M^{1,2}$ on the input to that step, and the observation is determined independently by a channel $C^{1,2}$ on the same input, i.e. *before* application of the markov.

FIGURE 1. Two successive steps $H^1$ and $H^2$ of a heterogeneous *HMM*.

We then have $\mathsf{avg}.(\llbracket C \rrbracket . \pi) = \pi$, because

$$(\mathsf{avg}.(\llbracket C \rrbracket . \pi))_x \;=\; (\mathsf{avg}.\llbracket \pi \triangleright C \rrbracket)_x \;=\; \Sigma(\pi \triangleright C)_{x,-} \;=\; \pi_x \quad .$$

In fact $\mathsf{avg}.\llbracket J \rrbracket$ for any $J$ in $\mathcal{X} \to \mathcal{Y}$ is $J$'s $\mathcal{X}$ marginal in $\mathbb{D}\mathcal{X}$.

## 3. Classical- vs. abstract *HMM*'s

3.1. **Classical *HMM*'s, and single *HMM*-steps as matrices.** Classically a Hidden Markov Model comprises a set $\mathcal{X}$ of states, a set $\mathcal{Y}$ of observations and two stochastic matrices $C, M$ that give resp. the *emission* probabilities $C_{x,y}$ that $x$ will emit observation $y$ and the *transition* probabilities $M_{x,x'}$ that $x$ will change to $x'$ [12]. Usually, the *homogeneous* case, computation evolves in (probabilistic) steps each determined by the same $C, M$, with each output state $x'$ becoming the following input $x$ and with the emissions $y$ accumulating. In our case however, *heterogeneous*, we can vary the matrices from step to step, each standing for various (different) program fragments.

We show two computations in Fig. 1. If $\pi$ is the distribution of incoming $x$, the distribution $\pi''$ of intermediate $x''$ is $\pi \cdot M^1$. The distribution of observations $y^1$ is $\pi \cdot C^1$. The second step's input $x''$ is the output of the first step.

A classical *HMM* hides all of three of $x, x'', x'$, but still the observations $y^{1,2}$ tell us something about each of them provided we know $\pi, M^{1,2}, C^{1,2}$. (This is analogous to knowing the source code of a program, but not being able to observe its variables as it executes.)

From now on we call the emission part of an *HMM* the *channel* and the transition part the *markov* (lower case).

**Definition 3.1** (Single *HMM*-step). Given channel $C \colon \mathcal{X} \to \mathcal{Y}$ and markov $M \colon \mathcal{X} \to \mathcal{X}$, define the *HMM*-matrix $(C;M)$ of type $\mathcal{X} \to \mathcal{Y} \times \mathcal{X}$ by

$$(C;M)_{x,y,x'} \;:=\; C_{x,y} \times M_{x,x'} \;.$$

---

[6]This $\mathsf{avg}$ is multiplication $\mu$ from the Giry monad: see §7.

```
// xs is initialised uniformly at random.
xs:= xs 1/2⊕ -xs
// What does an attacker guess for xs finally?
```

> The secret two-bit bit-string xs is set initially from $\{00, 01, 10, 11\}$ with equal probability $1/4$ for each; the following assignment either leaves xs unchanged (probability $1/2$) or bit-wise inverts both components.

FIGURE 2. Pure-markov *HMM* program

This (row-1-summing) matrix $(C;M)$ produces a joint distribution of type $\mathbb{D}(\mathcal{X} \times \mathcal{Y} \times \mathcal{X})$, as top-left in Fig. 6, once applied to a prior (Def. 2.2).

Note that in $(C;M)$ the probabilistic choices in $C$ (of $y$) and $M$ (of $x'$) are made independently; although indeed $(C;M)$ has the property that for each $x\colon\mathcal{X}$ the (remaining) joint distribution $(C;M)_{x,-,-}$ is independent in $y, x'$, this property is not preserved once steps are composed (§4).

3.2. **Abstract *HMM*'s represent classical *HMM*'s.** For abstract channels (§2.5) we focussed on the hyper of posteriors on the *input*; for *HMM*'s we focus on the hyper of posteriors on the *output*, because *HMM*'s are computations and so it is over their outputs we wish to reason. (The *prior* on the output would be our calculation from the input prior and the markov of what the output distribution would be, but *before* running the program and making observations in the type $\mathcal{Y}$.)

**Definition 3.2** (Matrix *HMM* denotes abstract *HMM*)**.** Let $H\colon \mathcal{X}{\rightarrow}\mathcal{Y}{\times}\mathcal{X}$ be an *HMM* presented as a matrix (stochastic in $y, x'$). Its denotation, of type $\mathbb{D}\mathcal{X}{\rightarrow}\mathbb{D}^2\mathcal{X}$, is called an *abstract HMM* and is defined $[\![H]\!].\pi := [\![J]\!]$, where $\pi\colon\mathbb{D}\mathcal{X}$ and the joint-distribution matrix $J\colon\mathcal{X}{\rightarrow}\mathcal{Y}$ is given by $J_{x',y} := \Sigma_x \pi_x H_{x,y,x'}$.

In §12 we discuss the (Dalenius) implications of having abstracted from the *HMM*'s input (with the $\Sigma_x$ just above) — it is not always appropriate.

3.3. **Special cases of *HMM*-steps: pure markovs.** Markovs are the special case of *HMM* where the channel-part effectively outputs nothing. If an *HMM*-step $(C;M)$ has for its channel $C$ an all-one column vector nc, where nc stands for "null channel". Then $\mathcal{Y}$ is a singleton and $J$ becomes a column vector: i.e. $J_{x'} = \Sigma_x \pi_x M_{x,x'}$, so that in fact $J$ is the usual matrix product $\pi{\cdot}M$.

Taking nc as the default channel gives $[\![:M]\!].\pi = [\![\text{nc}:M]\!].\pi = [\pi{\cdot}M]$, the point hyper on $\pi{\cdot}M$. This more general framework simply "wraps" a $[-]$ around the final distributions; but it's that wrapping that enables treating markovs and channels within the same type. A general $H$ is a markov just when $\Sigma_{x'} H_{x,y,x'}$ is nc.

```
// xs is initialised uniformly at random.
leak xs[0] 1/2⊕ xs[1]
```

> The value of either bit 0 or bit 1 of xs is revealed; the attacker learns that value,
> but does not know which bit it came from.
> What should he guess for xs after execution in this case?

FIGURE 3. Pure-channel *HMM* program.

Consider the program of Fig. 2 whose single variable is a two-bit string xs. We model it with $\mathcal{X}=\{00, 01, 10, 11\}$; prior $\pi\colon \mathbb{D}\mathcal{X}$ is uniform, and its markov $M$ is as just below:

$$
\begin{array}{c}
\begin{array}{cccc} 00 & 01 & 10 & 11 \end{array} \\
\begin{array}{c} 00: \\ 01: \\ 10: \\ 11: \end{array}
\begin{pmatrix}
1/2 & 0 & 0 & 1/2 \\
0 & 1/2 & 1/2 & 0 \\
0 & 1/2 & 1/2 & 0 \\
1/2 & 0 & 0 & 1/2
\end{pmatrix}
\end{array}
$$

The output distribution is of course $\pi'=\pi{\cdot}M=\pi$, and so the attacker's guess of the final state is optimally any of the four values in $\mathcal{X}$: they are equally good.

This system viewed as an abstract *HMM* would give output *hyper* $\Delta' = [\![:M]\!].\pi = [\pi]$, in fact the *point* hyper on $\pi$ indicating that the attacker is certain (point-probability 1) that the posterior distribution $\pi'$ on the final value of xs is equal to the prior $\pi$ in this case, i.e. it is still uniform.

3.4. **Special cases of *HMM*-steps: pure channels.** Channels are the special case where the input- and the output state are the same. If $(C;M)$ has markov $M$ as the identity $\mathsf{id}$, then it is a "pure channel" with output the same as its input. In that case Def. 3.2 gives $J_{x',y} = \sum_x \pi_x C_{x,y} \mathsf{id}_{x,x'} = (\pi{\rhd}C)_{x',y}$, and so $[\![C{:}\mathsf{id}]\!]$ from Def. 3.2 is just $[\![C]\!]$ from Def. 2.8.
  With $\mathsf{id}$ as the default markov, we have $[\![C{:}]\!] = [\![C]\!]$.
  Now consider Fig. 3 where some of xs is leaked, but xs itself is not changed. Thus our state $\mathcal{X}$ and prior $\pi$ are as before, the observation space is $\mathcal{Y}=\{0, 1\}$ and the channel $C$ representing this program is here at left:

$$
C = \begin{array}{c} \\ 00: \\ 01: \\ 10: \\ 11: \end{array}
\begin{array}{c}
\begin{array}{cc} 0 & 1 \end{array} \\
\begin{pmatrix}
1 & 0 \\
1/2 & 1/2 \\
1/2 & 1/2 \\
0 & 1
\end{pmatrix}
\end{array}
\qquad
J = \begin{array}{c} \\ 00: \\ 01: \\ 10: \\ 11: \end{array}
\begin{array}{c}
\begin{array}{cc} 0 & 1 \end{array} \\
\begin{pmatrix}
1/4 & 0 \\
1/8 & 1/8 \\
1/8 & 1/8 \\
0 & 1/4
\end{pmatrix}
\end{array}
$$

The joint distribution in $x', y$ is $J=\pi{\rhd}C$. The construction of Def. 2.7 gives us a hyper $\Delta'$ as

$$
\begin{array}{ll}
\text{inner distributions} & \text{outer distribution} \\
(1/2, 1/4, 1/4, 0) & @\,1/2 \\
(0, 1/4, 1/4, 1/2) & @\,1/2\ ,
\end{array}
\tag{3.1}
$$

where in general we write $z_1@p_1,\ z_2@p_2, \cdots$ for the discrete distribution that assigns probability $p_1$ to $z_1$ etc. In (3.1) the values $z_1, z_2$ are themselves (inner, posterior) distributions. This hyper shows that with probability $1/2$ an optimal attacker will guess 00 (because she

saw a 0 leaked, and deduces *a posteriori* that 00 now has the highest probability, twice either of the others); and with probability $1/2$ the attacker will guess 11 (because she saw a 1).

## 4. Beyond steps: *HMM* programming and sequential composition

### 4.1. Classical *HMM* composition: matrices.
Let $H^1, H^2 \colon \mathcal{X} \rightarrow \mathcal{Y} \times \mathcal{X}$ be two *HMM*'s. Their sequential composition $H = H^1; H^2$ describes the distribution on $x$, and $y_{1,2}$ together, and $x'$ as

$$(H^1; H^2)_{x,(y_1,y_2),x'} \quad := \quad \sum_{x''} H^1_{x,y_1,x''} H^2_{x'',y_2,x'} \ . \tag{4.1}$$

Note how the set of observables is now $\mathcal{Y} \times \mathcal{Y}$, compounding the observations $\mathcal{Y}$ from each component. (This is why infinite composition of *HMM*'s cannot easily be represented as a finite matrix.)

Remarkably, the action of *HMM*-composition on pure-markovs *HMM*'s is effectively their matrix multiplication, yet its action on pure channels is effectively their "parallel composition": thus a single general definition of composition specialises automatically to the two principal sub-cases, as we now show. First, we give the details for classical *HMM*'s; then Thm. 4.3 shows that the same holds for abstract *HMM*'s.

### 4.1.1. *Composition of pure markovs.*
The usual composition of Markov matrices $M^{1,2} \colon \mathcal{X} \rightarrow \mathcal{X}$ is via matrix multiplication $M^1 \cdot M^2$, and the result is of the same type $\mathcal{X} \rightarrow \mathcal{X}$. If we do it at the *HMM*-level, we find

$$
\begin{aligned}
&\quad (;M^1);(;M^2) \ _{x,(y_1,y_2),x'} \\
&= \quad \textstyle\sum_{x''} (;M^1)_{x,y_1,x''} (;M_2)_{x'',y_2,x'} \\
&= \quad \mathsf{nc}_{x,(y_1,y_2)} \textstyle\sum_{x''} M^1_{x,x''} M^2_{x'',x'} \qquad\qquad \text{``Recall from §3.3 that channel } \mathsf{nc} \text{ reveals nothing.''} \\
&= \quad (;\ M^1 \cdot M^2) \ _{x,(y_1,y_2),x'} \ ,
\end{aligned}
$$

so that indeed $(;M^1);(;M^2) = (;\ M^1 \cdot M^2)$.

### 4.1.2. *Composition of pure channels.*
Parallel composition of channels, which we write $C^1 \| C^2$, models applying *both* channels to the same input and observing both outputs. Thus

$$C^1 \| C^2 \quad _{x,(y_1,y_2)} \quad = \quad C^1_{x,y_1} \times C^2_{x,y_2} \ .$$

This is different from channel *cascading*, which applies the second channel $C^2$ to the observations of the first channel $C^1$ via matrix multiplication. A striking distinction is that the cascade of $C^1$ into $C^2$ releases no more information that $C^1$ alone (the Data-Processing Inequality [13]), whereas $C^1 \| C^2$ releases no *less* information that either of $C^{1,2}$ alone. In this latter case we find

$$
\begin{aligned}
&\quad (C^1;);(C^2;) \ _{x,(y_1,y_2),x'} \\
&= \quad \textstyle\sum_{x''} (C^1;)_{x,y_1,x''} (C^2;)_{x'',y_2,x'} \\
&= \quad \textstyle\sum_{x''} C^1_{x,y_1} \mathsf{id}_{x,x''} C^2_{x'',y_2} \mathsf{id}_{x'',x'} \qquad\qquad \text{``Recall from §3.4 markov } \mathsf{id} \text{ is the identity.''} \\
&= \quad C^1_{x,y_1} C^2_{x,y_2} \mathsf{id}_{x,x'} \\
&= \quad (C^1 \| C^2)_{x,(y_1,y_2)} \mathsf{id}_{x,x'} \\
&= \quad (C^1 \| C^2 \ ;) \ _{x,(y_1,y_2),x'} \ ,
\end{aligned}
$$

so that indeed again $(C^1;);(C^2;) = (C^1\|C^2\ ;\ )$.

4.1.3. *Pure channel followed by pure markov.* Finally, note that a general *HMM*-step (§3.1) is a pure channel followed by a pure markov. Let $\mathsf{nc}_{x,y_2}$ be (1 if $y_2=\hat{y}$ else 0) for some fixed $\hat{y}$ in $\mathcal{Y}$, and calculate

$$
\begin{array}{lll}
 & (C;);(;M)\quad_{x,(y_1,y_2),x'} & \\
= & \sum_{x''}(C;)_{x,y_1,x''}(;M)_{x'',y_2,x'} & \\
= & \sum_{x''} C_{x,y_1}\mathsf{id}_{x,x''}\mathsf{nc}_{x'',y_2}M_{x'',x'} & \\
= & C_{x,y_1}M_{x,x'}\mathsf{nc}_{x,y_2} & \text{``}\mathsf{id}_{x'',x'},\text{ 1-point rule''} \\
= & C_{x,y_1}M_{x,x'}\text{ if }y_2=\hat{y}\text{ else }0 & \text{``above''} \\
= & (C;M)_{x,(y_1,y_2),x'}\ , &
\end{array}
$$

so that $(C;)\,;\,(;M) = (C;M)$.

The reason that $(C;);(;M)$ and $(;M);(C;)$ differ in general is that in the (mathematical) definition of an *HMM*-step (e.g. Fig. 1) the emissions are determined by the input, initial state (rather than the output, final state). Had that original definition been the other way around, then we'd have had $(;M);(C;)$ as an *HMM*-step.

4.1.4. *Pure markov followed by pure channel.* This cannot, in general, be reduced to a single *HMM*-step. In $(;M);(C;)$ let both $C,M$ be the identity. Then the observations and final state will be perfectly correlated, something that is not possible for single *HMM*-step $(C';M')$.

4.2. **Abstract *HMM*'s: Kleisli composition.** Now we consider $h_1;h_2$ where $h_{1,2}$ are abstract *HMM*'s. (We use upper-case for matrices and lower-case for denotations.) Because the components' types $\mathbb{D}\mathcal{X}\to\mathbb{D}^2\mathcal{X}$ do not match directly, i.e. the co-domain $\mathbb{D}^2\mathcal{X}$ from the left is not the domain $\mathbb{D}\mathcal{X}$ required on the right, we use Kleisli composition for that.[7]

**Definition 4.1** (Push-forward of a function)**.** Given sets $\mathcal{Z},\mathcal{Z}'$ and function $f\colon\mathcal{Z}\to\mathcal{Z}'$, we write $\mathbb{D}f$ for the *push-forward* of $f$, a "lifted" function of type $\mathbb{D}\mathcal{Z}\to\mathbb{D}\mathcal{Z}'$ [14]. For $z'\colon\mathcal{Z}'$ and $\delta\colon\mathbb{D}\mathcal{Z}$ we have [8]

$$
\mathbb{D}f.\delta.z'\quad:=\quad\sum_{\substack{z:\,\mathcal{Z}\\ f.z=z'}}\delta.z\ .\qquad\text{[9]}
$$

**Definition 4.2** (Kleisli composition of abstract *HMM*'s)**.**
Given two abstract *HMM*'s $h_{1,2}\colon\mathbb{D}\mathcal{X}\to\mathbb{D}^2\mathcal{X}$, their Kleisli composition is defined

$$
(h_1;h_2).\pi\quad:=\quad\mathsf{avg}.(\mathbb{D}h_2.(h_1.\pi))
$$

for $\pi\colon\mathbb{D}\mathcal{X}$, where $\mathbb{D}h_2$ is as above the push-forward of $h_2$. Using functional composition, equivalently $h_1;h_2:=\mathsf{avg}\circ\mathbb{D}h_2\circ h_1$.

---

[7]This is the usual composition in a Kleisli category. See §7.

[8]Lifting, as in $\mathbb{D}f$, binds tightest: the conventional notation for $\mathbb{D}f.\delta.z'$ would be $(\mathbb{D}f)(\delta)(z')$, so that $(\mathbb{D}f)(\delta)\in\mathbb{D}\mathcal{Z}'$.

[9]$\mathbb{D}f$ is the action of functor $\mathbb{D}$ on arrow $f$: see §7.

That is, the lifting inherent in Kleisli-composition applies the right-hand abstract *HMM* $h_2$ to each inner (i.e. posterior) produced by the left-hand $h_1$ from prior $\pi$, preserving the way in which they are all combined together by the outer distribution. Then the intermediate result, of type $\mathbb{D}^3\mathcal{X}$, is averaged to bring it back to the required type $\mathbb{D}^2\mathcal{X}$.

4.3. **Proof that composition is faithfully denoted.** It is important (though unsurprising) for our interpretation that composition of *HMM*'s as matrices (4.1) is correctly mapped by $[\![\cdot]\!]$ to their Kleisli composition as abstract *HMM*'s (Def. 4.2). That is, we expect

**Theorem 4.3** (Composition faithfully denoted). *Let $H^{1,2}\colon \mathcal{X} \rightarrowtail \mathcal{Y} \times \mathcal{X}$ be HMM's as matrices. Then we have*

$$[\![H^1; H^2]\!] \quad = \quad [\![H^1]\!]; [\![H^2]\!] \ ,$$

*where (4.1) is used on the left and Def. 4.2 on the right.*

*Proof.* We reason as follows for any $\pi$.

$$
\begin{aligned}
&([\![H^1]\!]; [\![H^2]\!]).\pi \\
=\ &\mathsf{avg}.(\mathbb{D}[\![H^2]\!].([\![H^1]\!].\pi)) \\
=\ &\mathsf{avg}.(\mathbb{D}[\![H^2]\!].[\![[\![\pi]\!]H^1]\!]) \\
=\ &\mathsf{avg}.(\mathbb{D}[\![H^2]\!].(\textstyle\sum_{y^1} [([\![\pi]\!]H^1)_{y^1,-}])) \\
=\ &(\textstyle\sum_{y^1} [\![H^2]\!].([\![\pi]\!]H^1)_{y^1,-}) && \text{``}\mathbb{D}f \text{ distributes } f \text{ through inners.''} \\
=\ &(\textstyle\sum_{y^1} [\![([\![\pi]\!]H^1)_{y^1,-}]\!] H^2]\!]) \\
=\ &(\textstyle\sum_{y^1} (\textstyle\sum_{y^2} [([\![([\![\pi]\!]H^1)_{y^1,-}]\!] H^2)_{y^2,-}])) \\
=\ &(\textstyle\sum_{y^1,y^2} [([\![\pi]\!](H^1; H^2))_{(y^1,y^2),-}]) && \text{``Lem. 4.4''} \\
=\ &[\![H^1; H^2]\!].\pi \ ,
\end{aligned}
$$

as required.  □

**Lemma 4.4** (Double application of *HMM* matrix). *We have*

$$([\![\pi]\!](H^1; H^2))_{(y^1,y^2),-} \quad = \quad ([\![([\![\pi]\!]H^1)_{y^1,-}]\!]H^2)_{y^2,-}$$

*from this calculation for any $x'$ that*

$$
\begin{aligned}
&([\![\pi]\!](H^1; H^2))_{(y^1,y^2),x'} \\
=\ &(\textstyle\sum_x \pi_x (H^1; H^2)_{x,(y^1,y^2),x'}) \\
=\ &(\textstyle\sum_x \pi_x (\textstyle\sum_{x''} H^1_{x,y^1,x''} H^2_{x'',y^2,x'})) \\
=\ &(\textstyle\sum_{x,x''} \pi_x H^1_{x,y^1,x''} H^2_{x'',y^2,x'}) \\
=\ &(\textstyle\sum_{x''} (\textstyle\sum_x \pi_x H^1_{x,y^1,x''}) H^2_{x'',y^2,x'}) \\
=\ &(\textstyle\sum_{x''} ([\![\pi]\!]H^1)_{y^1,x''} H^2_{x'',y^2,x'}) \\
=\ &([\![([\![\pi]\!]H^1)_{y^1,-}]\!]H^2)_{y^2,x'} \ ,
\end{aligned}
$$

*as required.*

```
// xs is set uniformly at random.
leak xs[0] 1/2⊕ xs[1] ;
xs:= xs 1/2⊕ -xs
```

> The value of either bit 0 or bit 1 of `xs` is revealed; the attacker learns that value, but does not know which bit it is. Then `xs` is either unchanged or inverted, but the attacker does not know which.
>
> What's his best guess now for the final value of `xs`?

FIGURE 4. *HMM* program as sequential composition.

4.4. **Channel/markov together: two examples of composition.** For an example of sequential composition we return to `xs` and consider Fig. 4 where the state is both leaked *and* (possibly) changed. The final hyper $\Delta'$ in this case is obtained by applying the markov $M$ to the *inners* generated by $C$ in §3.3 while retaining their outers: that gives

$$
\begin{array}{ll}
(\quad 1/2 \times 1/2 + 1/2 \times 0, & (\quad 1/2 \times 0 + 1/2 \times 1/2, \\
1/2 \times 1/4 + 1/2 \times 1/4, & 1/2 \times 1/4 + 1/2 \times 1/4, \\
1/2 \times 1/4 + 1/2 \times 1/4, & 1/2 \times 1/4 + 1/2 \times 1/4, \\
1/2 \times 0 + 1/2 \times 1/2 \quad )@\,1/2 & 1/2 \times 1/2 + 1/2 \times 0 \quad )@\,1/2
\end{array}
$$

which is simplified first to this

$$
\begin{array}{ll}
(1/4, 1/4, 1/4, 1/4) & @\,1/2 \\
(1/4, 1/4, 1/4, 1/4) & @\,1/2
\end{array}
$$

and then, since the two inners are the same, as a hyper-distribution is collapsed to just the singleton hyper $[\pi]$, where we are using an explicit $(\times)$ for multiplication of specific numbers.

Thus the program of Fig. 4 reveals nothing about the final value of `xs` when the initial distribution was uniform. Informally we would explain this by noting that the information about `xs` released by the `leak` becomes "stale", irrelevant once we do not know whether `xs` has subsequently been inverted or not. (See §12 however for a discussion of why the initial value of `xs` might in some cases still be important.)

It would be wrong however to conclude, from $\Delta'=[\pi]$ in this specific case, that the program is secure for `xs` in general — for when the initial distribution is *not* uniform, the final value of `xs` *can* be less secure than the initial. This illustrates the danger in assuming something is uniformly distributed simply because we know nothing about it. (See §4.4.)

We now reconsider Fig. 4 but with a non-uniform prior, showing that indeed the conclusion that the program was (wrt. the final state) "leak free" is unjustified. In Fig. 5 the initial hyper is "skewed", i.e. it is not uniform over the whole type $XS$ of `xs`, but rather is concentrated on only three of its values:

$$
(0, 1/3, 1/3, 1/3) \qquad @\,1 \ , \tag{4.2}
$$

so that with certainty (@1) it is known that the initial distribution is $(0, 1/3, 1/3, 1/3)$. Via the first statement `leak xs[0]` $_{1/2}\oplus$ `xs[1]` an attacker will with probability $1/3$ (resp. $2/3$) observe 0 (resp. 1) and revise his belief of `xs`'s distribution as in the first (resp. second) row here:

$$
\begin{array}{ll}
(0, 1/2, 1/2, 0) & @\,1/3 \\
(0, 1/4, 1/4, 1/2) & @\,2/3
\end{array}
$$

```
// xs is set uniformly from {01,10,11}.
leak xs[0] 1/2⊕ xs[1] ;
xs:= xs 1/2⊕ -xs
```

> This system is as in Fig. 4 except that the prior initial distribution differs: at least one bit of xs is known to be 1.

FIGURE 5. Simple-channel program excluding xs=00 initially.

And after the second statement xs:= xs $_{1/2}\oplus$ -xs the hyper for the *current* (and final) distribution of xs will have become

$$
\begin{array}{ll}
(0, 1/2, 1/2, 0) & @\,1/3 \\
(1/4, 1/4, 1/4, 1/4) & @\,2/3 ,
\end{array}
\tag{4.3}
$$

where in the $1/3$-case he is better off finally than initially (since he knows xs cannot be 00 or 11), but in the other case he is worse off (since xs=00 has become possible). Thus if the attacker's choice is either to guess xs's initial value or to run the program and guess xs's final value, he can use these hypers to help make up his mind depending on his own criteria for the utility of his planned theft, that is the social context in which he is operating. Compare for example a thief's two alternatives for stealing a credit card: she might *"Steal it now, since the wallet is just sitting there."* or she might *"Steal it after the card is used at an ATM where she can see some digit of the PIN."* But in the second case there is a risk her victim will notice her, and choose a new PIN.

For example, the Shannon entropy of xs is initially $\lg(3){\sim}1.6$, but finally, it is conditionally $1/3{\times}1 + 2/3{\times}2 = 2/3{>}1.6$: if the attacker is using Shannon entropy to make his decision, he should act sooner rather than later.

On the other hand, the one-guess probability (Rényi min-entropy) of xs is initially $1/3$; and finally it is the same, at $1/3{\times}1/2 + 2/3{\times}1/4 = 1/3$. If the attacker is using this criterion, it does not matter when he acts.

In either case, the hypers (4.2) and (4.3) contain all the information necessary for his decision: the bit-values printed are themselves not important *for his decision*, which is why we can quotient our semantics by abstracting from them. (He does, however, need those values when he *makes* his attack if indeed he decides "later".)

These calculations are confirmed in the next section.

## 5. OVERVIEW OF HASKELL-MONADIC PROTOTYPE

A Haskell prototype of our hyper-based monadic model has been constructed for discrete, finite *HMM*'s, and it has been applied to our examples of Figs. 2–5 [15]. We give a brief summary here.

A discrete probability distribution on a set $\mathcal{X}$ is modelled as a monadic type Dist x that is effectively a list [(x,Rational)] of elements from $\mathcal{X}$ and their associated probabilities. The type of (discrete) hypers $\mathbb{D}^2\mathcal{X}$ is then Dist(Dist x).

A Markov "matrix" on $\mathcal{X}$ is of type x->Dist x, in fact encoding the matrix as a function from row-indices to distributions $\mathbb{D}\mathcal{X}$; a channel matrix is of type x->Dist y for any type $\mathcal{Y}$ of observations whatever.

The mini- programming language has two elementary statements: to use a markov `mm` we have an `update mm` that updates the state according to `mm`. Note that `mm` is a *Markov matrix*, but `update mm` is a *markov HMM* that is constructed from `mm`, i.e. implements it.

To use a channel `cm` we have a `reveal cm` that emits (e.g. prints) the channel's output wrt. the hidden state at that point: the state is not changed and, in particular, the output is not assigned to anything. It is merely observed. Both of these statements are of type `Dist x->Dist(Dist x)`, modelling our $\mathbb{D}\mathcal{X}\rightarrow\mathbb{D}^2\mathcal{X}$. In fact they are in $\mathbb{H}\mathcal{X}$, as Lemmas 8.1,8.3 show.

Sequential composition (;) of programs is the Kleisli composition `>=>` provided by Haskell's conventions for monads, in this case the monad `Dist`. Using that, and relying on §4.3 and §4.1.3, we can define an elementary *HMM*-step (§3.1) as

```
      hmmStep cm mm
    = reveal cm >=> update mm .
```

Thus `hmmStep` does not have to be primitive.

Our example space $\mathcal{X}$ is `(Bool,Bool)`, representing the bit-pair `xs`, and our two example (input) priors are

```
uniform      = [  ((False,False),1%4),
                  ((False,True),1%4),
                  ((True,False),1%4),
                  ((True,True),1%4)
                ]
```

*and from* §B

```
skewed       = [  ((False,False),0),
                  ((False,True),1%3),
                  ((True,False),1%3),
                  ((True,True),1%3)
                ]
```

Our example observation space $\mathcal{Y}$ is `Bool`.

With suitable definitions typed as above for channel `oneBit` that outputs one of `xs`'s two bits, uniformly at random, and `invert` that either inverts `xs` or does not, again uniformly at random, the four programs are then

```
      fig2 =       update invert
      fig3 =       reveal oneBit

      fig4 =       -- Uniform prior.
                   reveal oneBit
               >=> update invert

      fig5 =       -- Skewed prior.
                   reveal oneBit
               >=> update invert
```

The programs are run using the function

```
        runOn prior prog
    = pretty (prog prior)
```

where `pretty` is an output-formatting function that prints hypers in a readable way.

The results of running the programs are as follows, where the third column gives the outer probabilities of the resulting hyper, and the first two columns give the corresponding inner distributions. We print `True,False` as `1,0` respectively: (The prototype prints probabilities as fractions; but here they are printed as reals, for neatness.)

```
runOn uniform fig2 =  00  0.25  1.0      point hyper
                      01  0.25
                      10  0.25
                      11  0.25


runOn uniform fig3 =  01  0.25  0.5      Half the time...
                      10  0.25
                      11  0.5            11 is most likely, and

                      00  0.5   0.5      the other half it's 00.
                      01  0.25
                      10  0.25


runOn uniform fig4 =  00  0.25  1.0
                      01  0.25
                      10  0.25
                      11  0.25


runOn skewed fig5 =   00  0.25  0.67
                      01  0.25
                      10  0.25
                      11  0.25

                      01  0.5   0.33
                      10  0.5
```

The prototype contains also a `repeat` feature: for example `repeat 10 (reveal oneBit)` is a program that reveals a random bit of `xs` 10 times independently. (Such an iteration of parallel compositions is sometimes called "repeated independent runs.") With the uniform prior we would expect that the resulting hyper would have three inners: one of them, occurring with probability approximately $1/2$, would correspond to the case where the input bits of `xs` differed, in which case with overall probability $1023/1024$ there would be two different revelations among the 10 instances — thus showing that indeed the bits differed. But we would still have no (more) information about whether the input was `01` or `10`.

The remaining $1/1024$ would be split between two cases: bit `xs[0]` was revealed every time, or `xs[1]` was; and those outcomes would contribute to the other two inners.

Those other two inners would have probability approximately $1/4$ each, corresponding to input `00` or `11` where the two bits are the same. The program confirms this, giving

```
runOn uniform
(repeat 10 (reveal oneBit)) =
```

```
01  1/1026    513/2048   outer is about 1/4
10  1/1026
11  512/513   inner is "almost certainly 11"

01  1/2       511/1024   about 1/2
10  1/2

00  512/513   513/2048   about 1/4
10  1/1026
11  1/1026
```

(where this time we preserve the fractions). The small perturbations away from $1/4$ etc. reflect the small chance, mentioned above, that even when the inputs differ the random `oneBit` reveals the same bit 10 times in a row.

Finally, if we run the same program but with the final probabilistic inversion included, we get

```
runOn uniform
(repeat 10 (reveal oneBit)
>=> update invert) =

00  256/513   513/1024   two inners merged
01  1/1026
10  1/1026
11  256/513   about 1/2

01  1/2       511/1024   about 1/2
10  1/2
```

in which the two "bits equal" inners from just above have merged: although the probabilistic inversion preserves the information concerning whether the bits are equal, it conceals in the equals case whether they were both 00 or both 11.


## 6. The structure of hyper-space

Our hyper-space $\mathbb{D}^2\mathcal{X}$ has been synthesised by abstraction from the classical "matrix style" description of $HMM$'s. We now recall that there is a partial order ($\sqsubseteq$) of refinement on hypers, where for two hypers $\Delta_{S,I}:\mathbb{D}^2\mathcal{X}$ we say that $\Delta_S$ (a specification) is "refined by" $\Delta_I$ (implementation) when, in a sense we make precise below, the implementation $\Delta_I$ releases no more information than the specification $\Delta_S$ does [3–5,8]. That order lifts pointwise to $\mathbb{D}\mathcal{X}{\rightarrow}\mathbb{D}^2\mathcal{X}$, i.e. that $h_S{\sqsubseteq}h_I$ just when $h_S.\pi \sqsubseteq h_I.\pi$ for all $\pi:\mathbb{D}\mathcal{X}$, thus giving a new *refinement* order for (abstract) $HMM$'s. We write $\Delta_S{\sqsubseteq}\Delta_I$, and call it "uncertainty refinement" if we need to distinguish it from other kinds of refinement. Its ultimate antecedent is the lattice of information [16] — but it generalises those seminal ideas significantly.

**Definition 6.1** (Uncertainty refinement [3,5])**.** Let $\Delta_{S,I}:\mathbb{D}^2\mathcal{X}$ be two hypers on $\mathcal{X}$. We say that $\Delta_S$ is refined by $\Delta_I$ just when there is a distribution $\underline{\Delta}:\mathbb{D}^3\mathcal{X}$, that is a distribution of *hypers*, such that

$$\Delta_S \;=\; \mathsf{avg}.\underline{\Delta} \quad \text{and} \quad (\mathbb{D}\mathsf{avg}).\underline{\Delta} \;=\; \Delta_I \;.$$

Recall that $\mathbb{D}\mathsf{avg}$ is the push-forward of $\mathsf{avg}$ (Defs. 4.1,2.10).

The advantage of the abstract formulation in Def. 6.1 is that it is defined on hypers directly, and can be generalised to proper measures, thus extending discrete distributions [5].

But in the case (as here) where we remain discrete, there is an equivalent matrix-style characterisation:

**Lemma 6.2** (Refinement of joint-distributions [4, 8])**.** *Let $J_S: \mathcal{X} \twoheadrightarrow \mathcal{Y}_S$ and $J_I: \mathcal{X} \twoheadrightarrow \mathcal{Y}_I$ be joint-distribution matrices such that $[\![J_{S,I}]\!] = \Delta_{S,I}$ resp.* [10] *Then*

$$\Delta_S \sqsubseteq \Delta_I \quad \textit{iff} \quad J_S \cdot R = J_I \tag{6.1}$$

*for some stochastic refinement matrix $R: \mathcal{Y}_S \twoheadrightarrow \mathcal{Y}_I$. Note that the state-spaces of $\Delta_{S,I}$ are the same, but their observation spaces $\mathcal{Y}_{S,I}$ can differ.*

*Proof.* Illustrated in §C; sketch proof in §D.                                                                   □

With Lem. 6.2 the reflexivity and transitivity of relation ($\sqsubseteq$) is clear from elementary matrix properties. For anti-symmetry we refer to [6, Thm 6], whose supporting Lemma 1 there is adapted to suit our purposes here:

**Definition 6.3** (Expected value)**.** For distribution $\delta: \mathbb{D}\mathcal{Z}$ and function $f: \mathcal{Z} \to V$ for vector space $V$, the expected value of $f$ on $\delta$ is $\mathcal{E}_\delta f := \sum_{z: \mathcal{Z}} \delta_z \times f.z$, where $\sum$ and ($\times$) are taken in the vector space. [11]

We will be using $\mathcal{E}$ principally over hypers, i.e. the case $\mathcal{Z} = \mathbb{D}\mathcal{X}$ in the definition.

**Lemma 6.4** ((Strict) monotonicity)**.** *Given are two hypers $\Delta_{S,I}: \mathbb{D}^2\mathcal{X}$ and a strictly concave function $f: \mathbb{D}\mathcal{X} \to \mathbb{R}^{\geq}$.*
*If $\Delta_S \sqsubset \Delta_I$ then $\mathcal{E}_{\Delta_S} f < \mathcal{E}_{\Delta_I} f$. And if $f$ is (non-strictly) concave, then $\Delta_S \sqsubseteq \Delta_I$ implies $\mathcal{E}_{\Delta_S} f \leq \mathcal{E}_{\Delta_I} f$.*

*Proof.* Proved for abstract channels in [6, Lem 1]; the proof for hypers is essentially identical.
                                                                                                                    □

We now have antisymmetry, because $\Delta_S \sqsubseteq \Delta_I \sqsubseteq \Delta_S$ and $\Delta_S \neq \Delta_I$ implies $\Delta_S \sqsubset \Delta_I \sqsubset \Delta_S$ whence we have from Lem. 6.4 the contradiction $\mathcal{E}_{\Delta_S} f < \mathcal{E}_{\Delta_I} f < \mathcal{E}_{\Delta_S} f$ for any strictly concave $f: \mathbb{D}\mathcal{X} \to \mathbb{R}^{\geq}$ of our choice (for example Shannon entropy).

Hyper-space $\mathbb{D}^2\mathcal{X}$ also admits a metric, the Kantorovich metric [17] based on the Manhattan metric on $\mathbb{D}\mathcal{X}$ (§7). It is used for continuity properties (as we will see in §8.1), and is chosen because of its hierarchical properties, i.e. that the Kantorovich metric on say $\mathcal{X}$ induces a metric on $\mathbb{D}\mathcal{X}$ and $\mathbb{D}^2\mathcal{X}$ etc. [17].

## 7. Monads: Giry, Kleisli and Kantorovich

With $\mathbb{D}\mathcal{X} \to \mathbb{D}^2\mathcal{X}$ we have given a discrete model of abstract *HMM*'s, suitable for interpreting probabilistic sequential programs with hidden state, together with concrete programming examples (Figs. 2–4). We now provide a brief overview about how our setup embeds into structures based on a Giry monad.

The Giry monad over the category of Polish spaces and continuous functions comprises an endofunctor $\Pi$ and two natural transformations $\eta$ (unit) and $\mu$ (multiply) [2]; following [1] we take that as a basis for the denotation of computations. More precisely, we restrict

---

[10]Recall that the $\mathcal{X}$ here in type $\mathcal{X} \twoheadrightarrow \mathcal{Y}_S$ is the final-, not the initial state.

[11]More generally it is $\int f \, d\delta$ and requires measurability of $f$. One reason we do not use the standard notation $E(X)$ for the expected value of random variable $X$ is that the distribution over which $X$ is taken is implicit. In the calculations our *HMM*-semantics entails, we often need to make it explicit.

ourselves to the category *Comp* of compact metric spaces and continuous functions. We have been using $\mathbb{D}$ as a specialisation of $\Pi$ to this case. The object $\mathbb{D}\mathcal{S}$ is the set of Borel probability measures over the compact metric space $\mathcal{S}$ which is indeed a compact metric space [18, Thm 6.4]. To form a monad on *Comp*, we have provided the unit-function $[\cdot]$ specialising $\eta$ that makes a point measure, and multiply-function avg specialising $\mu$ that takes the average of a distribution (of distributions). Typically we have $[\cdot]_{\mathbb{D}\mathcal{X}} \in \mathbb{D}\mathcal{X} \to \mathbb{D}^2\mathcal{X}$ and $\mathsf{avg}_{\mathcal{X}} \in \mathbb{D}^2\mathcal{X} \to \mathbb{D}\mathcal{X}$ where the subscripts are left implicit when they are clear from the context.

From Giry's construction, the arrows $[\cdot]_{\mathcal{S}}$ and $\mathsf{avg}_{\mathcal{S}}$ are continuous with respect to the weak topology on $\mathbb{D}\mathcal{S}$ but, in this paper, we are dealing with compact metric spaces. Fortunately, for compact metric spaces, the Kantorovich distance metrizes the weak topology [19]. This implies that the triple $(\mathbb{D}, [\cdot], \mathsf{avg})$ is indeed a monad on the category *Comp*.

Monadic constructions based on the Kantorovich metric are not new. In [17], Van Breugel construct monads on the category *Comp\** of compact metric spaces and 1-Lipschitz functions. His functor $\mathcal{B}$ coincides with our $\mathbb{D}$ on objects and he shows that $\mathcal{B}f$, $[\cdot]_{\mathcal{S}}$ and $\mathsf{avg}_{\mathcal{S}}$ are 1-Lipschitz; whenever $f$ is 1-Lipschitz and $\mathcal{S}$ is a compact metric space. Thus $(\mathcal{B}, [\cdot], \mathsf{avg})$ is a monad on the category *Comp\** [12]. In that work, the metric is crucial since the notion of 1-Lipschitzness is not a topological property. In fact, Van Breugel shows that the Kantorovich metric is *the right* metric to construct probabilistic monads out of metric spaces. Such a construction does not necessarily work with other metrics that metrizes the weak topology (e.g. Prohorov metric which is equivalent to the Kantorovich metric from compact spaces). However, arrows in the category *Comp\** are insufficient to denote probabilistic programs with hidden states because $[\![C]\!]$ is not necessarily 1-Lipschitz for some channel matrix $C$. This drives our choice of the Giry monad $(\mathbb{D}, [\cdot], \mathsf{avg})$ on the category *Comp* of compact metric spaces and continuous functions.

Our construction of the Kantorovich metric begins with a finite set $\mathcal{X}$ endowed with the discrete metric $d_1$ (i.e. $d_1(x, x') := 0$ if $x = x'$ else 1). This is trivially a compact metric space. The space $\mathbb{D}\mathcal{X}$ of discrete distributions on $\mathcal{X}$ is endowed with the Kantorovich metric based on $d_1$ which coincides with the total variation metric on $\mathbb{D}\mathcal{X}$. At this level, the Kantorovich metric reduces to $d_K(\delta^1, \delta^2) = \frac{1}{2}\sum_x |\delta_x^1 - \delta_x^2|$. At the next level, our hyper-space $\mathbb{D}^2\mathcal{X}$ has the Borel algebra generated by $\mathbb{D}$ from $\mathbb{D}\mathcal{X}$, which is in turn determined by the Kantorovich metric derived from $d_K$ which we will also denote by $d_K$. These metrics are distinguished in terms of the arguments they are applied on, i.e. $d_K(\delta^1, \delta^2)$ is the Kantorovich metric on $\mathbb{D}\mathcal{X}$ while $d_K(\Delta^1, \Delta^2)$ is the Kantorovich metric on $\mathbb{D}^2\mathcal{X}$.

## 8. Characteristics of $\mathbb{H}\mathcal{X}$, the abstract *HMM*'s

### 8.1. Continuity and super-linearity.
The semantic function $[\![\cdot]\!]$ (Def. 3.2) takes *HMM* matrices in $\mathcal{X} \to \mathcal{Y} \times \mathcal{X}$ to functions in $\mathbb{D}\mathcal{X} \to \mathbb{D}^2\mathcal{X}$; but not all of those functions are denotations $[\![H]\!]$ of some $H$. We now describe two important characteristics satisfied by $[\![H]\!]$ as $H$ ranges over *HMM*'s: they are continuity and super-linearity. We will define abstract *HMM*'s $\mathbb{H}\mathcal{X}$ to be just the functions satisfying those conditions.

Sometimes called "healthiness conditions", these are essentially technical results giving general properties that are to hold for the denotation of any program. They are used for the

---

[12]Note that Van Breugel's construction is more general than this since he also considers complete metric spaces and tight probability measures over them.

proof of other, more specific properties of programs. Here we use them to define a subset $\mathbb{H}\mathcal{X}$ of $\mathbb{D}\mathcal{X}{\to}\mathbb{D}^2\mathcal{X}$, and we prove that $[\![H]\!]\in\mathbb{H}\mathcal{X}$ for all classical HMM's $H$.

Our first condition concerns continuity wrt. the Kantorovich metrics on $\mathbb{D}\mathcal{X}$ and $\mathbb{D}^2\mathcal{X}$.

**Lemma 8.1** (Denotations of *HMM*'s are continuous). *For all $H{:}\mathcal{X}{\to}\mathcal{Y}{\times}\mathcal{X}$ we have that $[\![H]\!]$ is a continuous function in $\mathbb{D}\mathcal{X}{\to}\mathbb{D}^2\mathcal{X}$ wrt. the Kantorovich metrics.*

*Proof.* We recall from Def. 3.2 that $[\![H]\!].\pi = [\![\pi{\triangleright}J]\!]$ where $J_{x',y} = \sum_x H_{x,y,x'}$. We consider this as the composition of the two functions $\pi{\mapsto}(\pi{\triangleright}J)$ and $[\![\cdot]\!]$, with the metric on the intermediate space $\mathcal{X}{\to}\mathcal{Y}$ (of matrices) being the Kantorovich metric $d_K$ on $\mathbb{D}(\mathcal{X}{\times}\mathcal{Y})$ which is $d_K(J^1, J^2) = \frac{1}{2}\sum_{(x,y)}|J^1_{x,y} - J^2_{x,y}|$. Furthermore, since $\pi{\mapsto}(\pi{\triangleright}J)$ comprises only elementary arithmetic operations, and $d_K$ is topologically equivalent to the Euclidean distance, the continuity is clear. Thus we concentrate on the continuity of $[\![\cdot]\!]$ at an arbitrary joint distribution $J{:}\mathbb{D}(\mathcal{X}{\times}\mathcal{Y})$.

Let $\varepsilon{>}0$. We denote $[\![J]\!] = \Delta$ and let $J'{:}\mathbb{D}(\mathcal{X}{\times}\mathcal{Y})$ with $[\![J']\!] = \Delta'$. Since $J$ and $J'$ are matrices, we can write $\Delta = \sum_y a_y[\delta^y]$ and $\Delta' = \sum_y a'_y[\delta'^y]$ with $\delta^y, \delta'^y{:}\mathbb{D}\mathcal{X}$. These are sums over the full set $\mathcal{Y}$ so if $a_y = 0$ (resp. $a'_y = 0$) then we define $\delta^y{:=}\delta'^y$ (resp. $\delta'^y{:=}\delta^y$). Let us define $\Delta'' = \sum_y a_y[\delta'^y]$ which combines the coefficients of $\Delta$ with the inners of $\Delta'$. The triangular inequality tells us that

$$d_K(\Delta, \Delta') \leq d_K(\Delta, \Delta'') + d_K(\Delta'', \Delta') \ .$$

On the one hand,

$$
\begin{array}{lll}
& d_K(\Delta'', \Delta') & \\
= & \frac{1}{2}\sum_y |a_y - a'_y| & \text{``}\Delta'' \text{ and } \Delta' \text{ have finite supports''} \\
= & \frac{1}{2}\sum_y |\sum_x J_{x,y} - \sum_x J'_{x,y}| & \text{``Defn. } a_y \text{ and } a'_y\text{''} \\
\leq & \frac{1}{2}\sum_{y,x} |J_{x,y} - J'_{x,y}| & \text{``}|\sum_x f.x| \leq \sum_x |f.x|\text{''} \\
= & d_K(J, J') & \text{``Defn. } d_K(J, J')\text{''}
\end{array}
$$

On the other hand, for every $y$, the function which maps $J$ to $\delta_y$ is continuous at $J$ since it is a composition of a $y$-projection and normalisation. Therefore, there exists $\alpha_y{>}0$ such that for every $J'{:}\mathbb{D}(\mathcal{X}{\times}\mathcal{Y})$ with $d_K(J, J'){<}\alpha_y$, we have $d_K(\delta^y, \delta'^y){<}\frac{\varepsilon}{2}$. But we have

$$
\begin{array}{lll}
& d_K(\Delta, \Delta'') & \\
\leq & \sum_y a_y d_K(\delta^y, \delta'^y) & \text{``Kantorovich-Rubeinstein Theorem [19, Pg. 8]''} \\
\leq & \max_y d_K(\delta^y, \delta'^y) & \text{``}\sum_y a_y = 1 \text{ and } a_y{\geq}0 \text{ for all } y\text{''}
\end{array}
$$

Therefore, we choose $\beta = \min(\min_y \alpha_y, \frac{\varepsilon}{2})$ and for every $J'$ such that $d_K(J, J') < \beta$, we have

$$d_K(\Delta, \Delta'){<}\frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon \ .$$

Hence, $[\![\cdot]\!]$ is continuous at $J$. □

Our second condition concerns linear combinations.

**Definition 8.2** (Weighted sum). For $\delta_{1,2}{:}\mathbb{D}\mathcal{X}$ we write $\delta_1{}_p{+}\delta_2$ for the weighted sum of the two distributions, so that $(\delta_p{+}\delta')_x = p\delta_x + (1{-}p)\delta'_x$. Note that $\delta_p{+}\delta'$ defined here and the $\delta_p{\oplus}\delta'$ of Def. 2.10 differ: the former is a single distribution made from $p$-merging $\delta, \delta'$; the latter is a hyper whose support is just the two elements $\delta, \delta'$.

**Lemma 8.3** (Denotations of *HMM*'s are super-linear)**.** *For all* $H\colon \mathcal{X} \twoheadrightarrow \mathcal{Y} \times \mathcal{X}$ *we have*

$$[\![H]\!].\pi_1 \ _p+ \ [\![H]\!].\pi_2 \quad \sqsubseteq \quad [\![H]\!].(\pi_1 \ _p+ \ \pi_2) \ , \tag{8.1}$$

*where* $(\sqsubseteq)$ *is refinement as defined in Def. 6.1.*[13]

*Proof.* Take any reduced $J^{1,2}\colon \mathcal{X} \twoheadrightarrow \mathcal{Y}$, and argue first that for any $0 \le p \le 1$ we have

$$[\![J^1]\!] \ _p+ \ [\![J^2]\!] \quad \sqsubseteq \quad [\![J^1 \ _p+ \ J^2]\!] \ , \tag{8.2}$$

since the horizontal concatenation $J$ of the two (scaled) matrices $p \times J^1$ and $(1{-}p) \times J^2$ satisfies $[\![J]\!] = [\![J^1]\!]_p + [\![J^2]\!]$, and $J$ itself is refined to $J_1 \ _p+ \ J_2$ (in the sense of Lem. 6.2) by the refinement matrix

$$R \quad := \quad \begin{pmatrix} 1 & 0 & \cdots \\ 0 & 1 & \cdots \\ \vdots & \vdots & \vdots \\ 1 & 0 & \cdots \\ 0 & 1 & \cdots \\ \vdots & \vdots & \vdots \end{pmatrix} \begin{array}{l} \longleftarrow \text{ corresp. to first col. of } J^1 \\ \longleftarrow \text{ corresp. to snd. col. of } J^1 \\ \\ \overset{\leftharpoonup}{\longleftarrow} \text{ corresp. to first col. of } J^2 \\ \longleftarrow \text{ corresp. to snd. col. of } J^2 \\ \\ \end{array}$$

that simply sums corresponding columns. Now we observe that

$$\begin{array}{lll}
& [\![H]\!](\pi_1 \ _p+ \ \pi_2) & \\
= & [\![(\pi_1 \ _p+ \ \pi_2) \triangleright H)]\!] & \\
= & [\![\pi_1 \triangleright H \ _p+ \ \pi_2 \triangleright H]\!] & \\
\sqsupseteq & [\![\pi_1 \triangleright H]\!] \ _p+ \ [\![\pi_2 \triangleright H]\!] & \text{``(8.2) just above''} \\
= & [\![H]\!].\pi_1 \ _p+ \ [\![H]\!].\pi_2 \ , &
\end{array}$$

as required.        $\square$

Motivated by those two lemmas, we now define

**Definition 8.4** (The space $\mathbb{H}\mathcal{X}$ of abstract *HMM*'s)**.** We write $\mathbb{H}\mathcal{X}$ for those $h$ in $\mathbb{D}\mathcal{X} \to \mathbb{D}^2\mathcal{X}$ satisfying Lemmas 8.1,8.3, i.e. that are Kantorovich-continuous and super-linear.

Thus our two lemmas above establish that $[\![H]\!] \in \mathbb{H}\mathcal{X}$ for any classical *HMM* $H$.

Since we will therefore be restricting our denotations to $\mathbb{H}\mathcal{X}$, a subset of the arrows in the category *Comp*, we expect $\mathbb{H}\mathcal{X}$ to be closed under composition.

**Lemma 8.5** (Abstract *HMM*'s closed under composition)**.** *For any two* $h_{1,2}\colon \mathbb{H}\mathcal{X}$ *we have* $h_1; h_2 \in \mathbb{H}\mathcal{X}$ *as well, where* (;) *is as in Def. 4.2.*

*Proof.* Although a direct proof is possible, the result is much easier once we have introduced "uncertainty" transformers (§9), because it is then a consequence of Thm. 9.8 and in particular its Cor. 9.9, which depends crucially on the dual view we develop in §9.      $\square$

It is shown in §E.2 that composition in $\mathbb{H}\mathcal{X}$ is monotonic with respect to the refinement order $(\sqsubseteq)$. This completes our construction of our forward, abstract semantics for *HMM*'s. We now propose a dual view.

---

[13]Super-linearity can also be seen as a form of monotonicity. See §E.1.

## 9. A dual view: uncertainty measures, and their transformers

### 9.1. Uncertainty measures, and their relation to refinement.

"Uncertainty measures" generalise the diversity of entropy measures (including e.g. Shannon), the functions from distributions to reals that measure increasing disorder.

**Definition 9.1** (Uncertainty measure). An *uncertainty measure* over $\mathcal{X}$ is a Kantorovich-continuous- and concave function in $\mathbb{D}\mathcal{X} \to \mathbb{R}^{\geq}$, i.e. one taking distributions (on $\mathcal{X}$ in this case) to non-negative reals. It is intended that a distribution's greater uncertainty indicates more resilience (less vulnerability) to the distribution's being exploited by an adversary. [14]

We write $\mathbb{U}\mathcal{X}$ for the uncertainty measures over $\mathcal{X}$, and call them "*UM*'s" in the text for brevity.

A typical example of a *UM* applied to a hyper is as follows. Given prior $\pi : \mathbb{D}\mathcal{X}$ and channel $C : \mathcal{X} \rightarrow \mathcal{Y}$, the resulting hyper is $\Delta := [\![\pi \triangleright C]\!]$ and the "conditional $u$ uncertainty" of that (compare conditional Shannon entropy) would be $\mathcal{E}_\Delta u$. We write $\mathcal{E}_\Delta u$ because it makes explicit that the conditional uncertainty is an expected value and, as such, we can calculate with it. (More conventional notations such as $H(y|x)$ –in the Shannon case– make those calculations more difficult.) This could be compared to the uncertainty $u.\pi$ of the prior, to give a "$u$-leakage" of the channel on that prior.

There is a compelling connection between *UM*'s (Def. 9.1) and refinement (Def. 6.1, Lem. 6.2): we have

**Lemma 9.2** (Soundness and completeness of uncertainty measures [6]). *For any hypers* $\Delta_{1,2} : \mathbb{D}^2\mathcal{X}$ *we have*

$$\Delta_1 \sqsubseteq \Delta_2 \quad \text{iff} \quad \mathcal{E}_{\Delta_1} u \leq \mathcal{E}_{\Delta_2} u \quad \text{for all } u : \mathbb{U}\mathcal{X}.$$

We regard "only if" as *soundness* in the sense that if we have a witness to the refinement relation $\Delta_1 \sqsubseteq \Delta_2$, i.e. either $\underline{\Delta}$ (Def. 6.1) or $R$ (Lem. 6.2), then no *UM* can show $\Delta_2$ to be less uncertain than $\Delta_1$. It is related to the *Data-Processing Inequality*, as explained in [6].

We regard "if" as *completeness* in the sense that if refinement fails, that is if $\Delta_1 \not\sqsubseteq \Delta_2$, then there is a *UM* demonstrating the failure [4–6,8].

In §E.5 is background on the proof of Lem. 9.2, whose completeness part was originally called "Coriaceous" because it was hard to prove [8].

### 9.2. Abstract *HMM*'s to *UM*-transformers.

In §3 we introduced a "forward" denotational view of *HMM*'s that takes initial distributions to final hypers. Here we take the dual view, where an *HMM* takes a "post-uncertainty" to a "pre-uncertainty".

**Definition 9.3** (Uncertainty-measure transformers). Take $h : \mathbb{H}\mathcal{X}$ and $u : \mathbb{U}\mathcal{X}$. Define the *uncertainty transformer* wp.$h$ of type $\mathbb{U}\mathcal{X} \to \mathbb{U}\mathcal{X}$ so that for any $u : \mathbb{U}\mathcal{X}$ and $\pi : \mathbb{D}\mathcal{X}$ we have

$$\text{wp}.h.u.\pi \quad := \quad \mathcal{E}_{h.\pi} u \; ,$$

where on the right we are taking the expected value of $u$ on the hyper $h.\pi$. (Because $u$ is continuous, it is measurable.)

---

[14]Smith's "vulnerability measure" based on Bayes Risk [9] is an uncertainty measure except that it goes in the opposite direction.

By analogy with weakest preconditions for ordinary sequential programming [20], a *UM*-transformer wp.$h$ takes a *UM* to be applied *after* $h$ and produces a *UM* that equivalently can be applied *before* $h$. (Compare also [21–23] for probabilistic/demonic sequential programs.) The idea (and utility) is in goal-directed reasoning: if one knows the program, and knows also the uncertainty that it must achieve, with the uncertainty transformer one determines the minimum uncdertainty that is necessary *before* the program begins its execution.

In Lem. 9.4 we show well definedness of Def. 9.3, that is that wp.$h.u$ is indeed in $\mathbb{U}\mathcal{X}$.

**Lemma 9.4** (Well-definedness of Def. 9.3)**.** *If $h\colon\mathbb{H}\mathcal{X}$ is an abstract HMM and $u\colon\mathbb{U}\mathcal{X}$ is a UM, then wp.$h.u$ is in $\mathbb{U}\mathcal{X}$.*

*Proof.* See §F.                                                                                       $\square$

### 9.3. Characteristic properties of wp.$h$.

For $h\colon\mathbb{H}\mathcal{X}$ the *UM*-transformer wp.$h$ has a number of characteristic properties.

**Lemma 9.5** (wp.$h$ is linear and total)**.** *For every $h\colon\mathbb{H}\mathcal{X}$ and $t=$ wp.$h$ we have that $t$ is:*
(1) linear *so that for $a_{1,2}\colon\mathbb{R}^{\geq}$ and $u_{1,2}\colon\mathbb{U}\mathcal{X}$ we have*

$$t.(a_1 u_1 + a_2 u_2) \quad = \quad a_1 t.u_1 + a_2 t.u_2 \; ;$$

(2) monotonic, *so that $t.u_1.\delta \geq t.u_2.\delta$ for every $u_1 \geq u_2$ with $u_{1,2}\colon\mathbb{U}\mathcal{X}$ and $\delta\colon\mathbb{D}\mathcal{X}$, where we lift $(\geq)$ pointwise; and*
(3) total, *so that $t.\mathbf{1}{=}\mathbf{1}$ where $\mathbf{1}.\delta{:=}1$ for all $\delta\colon\mathbb{D}\mathcal{X}$.*

*Proof.* These properties are immediate from Def. 9.3.                                       $\square$

A further property of *UM*-transformers is that they are 1-Lipshitz in a certain sense:

**Lemma 9.6** (wp.$h$ is 1-Lipschitz)**.** *Take $h\colon\mathbb{H}\mathcal{X}$ and define $t{:=}$ wp.$h$. Let $|\cdot|$ be absolute value. Then $t$ is 1-Lipschitz in the sense that*

$$\sup_{\delta\colon\mathbb{D}\mathcal{X}} |t.u_1.\delta - t.u_2.\delta| \quad \leq \quad \sup_{\delta\colon\mathbb{D}\mathcal{X}} |u_1.\delta - u_2.\delta| \; .$$

*Proof.* Consider arbitrary $u_{1,2}\colon\mathbb{U}\mathcal{X}$. We reason

$$
\begin{array}{lll}
& \sup_{\delta\colon\mathbb{D}\mathcal{X}} |\text{wp}.h.u_1.\delta - \text{wp}.h.u_2.\delta| & \\
= & \sup_{\delta\colon\mathbb{D}\mathcal{X}} |\mathcal{E}_{h.\delta}\, u_1 - \mathcal{E}_{h.\delta}\, u_2| & \\
\leq & \sup_{\delta\colon\mathbb{D}\mathcal{X}}\ \mathcal{E}_{h.\delta}\, |u_1 - u_2| & \text{``property of } |\cdot|\text{''} \\
\leq & \sup_{\delta\colon\mathbb{D}\mathcal{X}}\quad \mathcal{E}_{h.\delta}\, (\sup_{\delta'\colon\mathbb{D}\mathcal{X}} |u_1.\delta' - u_2.\delta'|) & \text{``}\mathcal{E}\text{ monotonic''} \\
= & \sup_{\delta\colon\mathbb{D}\mathcal{X}} |u_1.\delta - u_2.\delta| \; . & \text{``}\mathcal{E}_{h.\delta}\, \mathbf{1} = 1 \text{ and rename } \delta' \text{ to } \delta\text{''}
\end{array}
$$

$\square$

Motivated by those lemmas, we define uncertainty transformers to be exactly the functions in $\mathbb{U}\mathcal{X}{\to}\mathbb{U}\mathcal{X}$ that satisfy the properties listed.

**Definition 9.7** (The uncertainty transformers $\mathbb{T}\mathcal{X}$)**.** The uncertainty transformers $\mathbb{T}\mathcal{X}$ are the functions in $\mathbb{U}\mathcal{X}{\to}\mathbb{U}\mathcal{X}$ that satisfy Lems. 9.5,9.6.

We note that transformers $\mathbb{T}\mathcal{X}$ are closed under composition. In §E.3 we show that refinement for $\mathbb{T}\mathcal{X}$ is pointwise $(\leq)$.

9.4. **$UM$-transformers back to abstract $HMM$'s.** The function wp.$(\cdot)$ has been shown to be of type $\mathbb{H}\mathcal{X}{\to}\mathbb{T}\mathcal{X}$. Here we show that this correspondence is exact, i.e. that for every $t{:}\,\mathbb{T}\mathcal{X}$ there is an $h{:}\,\mathbb{H}\mathcal{X}$ such that $t{=}\mathrm{wp}.h$ and, moreover, that the $h$ is unique.

The following theorem thus establishes the exact correspondence between $\mathbb{H}\mathcal{X}$ and $\mathbb{T}\mathcal{X}$, giving an analogue for hidden-state probabilistic programs to the well-known correspondence between demonic relations and conjunctive predicate transformers [20] that the former correspond exactly to those functions from predicates to predicates that distribution conjunction. (A further example, generalising that, is the correspondence between demonic/probabilistic programs and super-linear expectation transformers [22, 24].)

**Theorem 9.8** (Characterisation of transformers). *For any $t{:}\,\mathbb{T}\mathcal{X}$ there is a unique $h{:}\,\mathbb{H}\mathcal{X}$ such that $t{=}wp.h$.*

*Proof.* Let $t{:}\,\mathbb{T}$. The construction of $h$ starts by showing that the transformer $t$ can be extended into a linear function over the space $\mathbb{C}\mathcal{X}$ of continuous functions from $\mathbb{D}\mathcal{X}$ to itself. This extension is executed in two phases. Firstly, we show that the set $\mathbb{U}\mathcal{X}$ of continuous, concave and non-negative function over $\mathbb{D}\mathcal{X}$ generates a sub-vector space of $\mathbb{C}\mathcal{X}$. Thus, the first stage is an algebraic extension of $t$ to that generated sub-space. This extension is necessarily unique by linearity (Lem. 9.5). The second stage is a topological extension where $\mathbb{C}\mathcal{X}$ is endowed with the norm uniform making it a Banach space. In fact, we show that the sub-vector space generated by $\mathbb{U}\mathcal{X}$ is a dense sub-algebra of $\mathbb{C}\mathcal{X}$ using the Stone-Weierstrass Theorem [25, Thm. 5]. Thus the second stage of the extension is also unique by continuity. Full technical details of these extensions are given in §G.

We now have a unique continuous linear function $\hat{t}$ from $\mathbb{C}\mathcal{X}$ to itself which coincides with $t$ on $\mathbb{U}\mathcal{X}$. We shall construct an $h$ such that wp.$h = t$.

Fix $\delta{:}\,\mathbb{D}\mathcal{X}$. The function
$$f \quad\mapsto\quad \tilde{t}.f.\delta$$
maps each continuous function $f{:}\,\mathbb{C}\mathcal{X}$ to $\tilde{t}.f.\delta$ is a positive linear functional on $\mathbb{C}\mathcal{X}$; moreover $\tilde{t}.\mathbf{1}.\delta = t.\mathbf{1}.\delta = 1$, thus $\mathbf{1} \mapsto 1$. Therefore, the Riesz Representation Theorem for linear functionals [18, Ch. 2 Thm. 5.8] implies that there exists a unique Borel probability measure $\Delta_\delta$ on $\mathbb{D}\mathcal{X}$ such that $\tilde{t}.f.\delta = \mathcal{E}_{\Delta_\delta}\, f$, for every $f{:}\,\mathbb{C}\mathcal{X}$.

Define $h.\delta := \Delta_\delta$ for each $\delta{:}\,\mathbb{D}\mathcal{X}$. We now check that $h$ has the required properties demanded by Lemmas 8.1, 8.3.

Continuity: For the continuity assumption in Lem. 8.1, we let $\delta_n$ be a sequence of distributions in $\mathbb{D}\mathcal{X}$ converging to $\delta{:}\,\mathbb{D}\mathcal{X}$ with respect to the Kantorovich metric on $\mathbb{D}\mathcal{X}$. It suffices to show that the limit of $d_K(h.\delta_n, h.\delta)$ is 0, as $n$ goes to infinity. Since $\mathbb{D}\mathcal{X}$ is compact, the Kantorivich metric metrizes the weak topology and it suffices to show that $h.\delta_n$ converges weakly to the Borel measure $h.\delta$. Let $f{:}\,\mathbb{C}\mathcal{X}$, we have
$$\mathcal{E}_{h.\delta_n}\, f = \tilde{t}.f.\delta_n$$
Since $\tilde{t}.f$ is also continuous, the sequence $\tilde{t}.f.\delta_n$ converges to $\tilde{t}.f.\delta$ and thus, $h.\delta_n$ coverges weakly to $h.\delta$. [15]

Super-linear: For the super-linearity assumption in Lem. 8.3, suppose that $t{=}\mathrm{wp}.h$ is in $\mathbb{T}\mathcal{X}$ and take arbitrary $\delta_{1,2}{:}\,\mathbb{D}\mathcal{X}$. Then we reason

---

[15]This proof crucially depends on the compactness of $\mathbb{D}\mathcal{X}$. For Polish spaces, we can achieve the same result but using a more general result by Rao [26, Thm. 3.1].

$$h.\delta_{1p} + h.\delta_2 \ \sqsubseteq\ h.(\delta_{1p} + \delta_2)$$

if $\quad \mathcal{E}_{(h.\delta_{1p}+h.\delta_2)}\, u\ \leq\ \mathcal{E}_{h.(\delta_{1p}+\delta_2)}\, u$ $\qquad\qquad$ "for all $u\colon \mathbb{U}\mathcal{X}$ Lem. 9.2 Coriaceous"

if $\quad \mathrm{wp}.h.u.\delta_{1p} + \mathrm{wp}.h.u.\delta_2 \leq\ \mathrm{wp}.h.u.(\delta_{1p}+\delta_2)$ $\qquad\qquad$ "Defn. wp.()"

if $\quad \mathrm{wp}.h.u \in \mathbb{U}\mathcal{X}\ ,$ $\qquad\qquad$ "Defn. $\mathbb{U}\mathcal{X}$"

which was our assumption. $\hfill\square$

With these characterisations, we now can prove two technical facts. In the discrete case (as earlier) they seem self-evident. In the more general setting, however, the work is mainly in ensuring well definedness (e.g. that only measurable functions are integrated, etc.) The first establishes the usual connection between composition, this time between the forward- and backward semantics; the second confirms that $\mathbb{H}\mathcal{X}$ is closed under composition (i.e. preserves continuity and super-linearity, as claimed in Lem. 8.5).

**Corollary 9.9** (Transformer composition). *For any $h_{1,2}\colon \mathbb{H}\mathcal{X}$ we have that also $h_1; h_2 \in \mathbb{H}\mathcal{X}$, and furthermore that $wp.(h_1; h_2) = wp.h_1 \circ wp.h_2$.*

*Proof.* Direct calculation shows that $\mathrm{wp}.(h_1; h_2) = \mathrm{wp}.h_1 \circ \mathrm{wp}.h_2$, although the working is intricate in the general (Giry) case. Well definedness of $h_1; h_2$ itself uses the simpler properties of (functional) composition on the transformer side. See §H. $\hfill\square$

Also, transformer composition respects refinement (§E.4).

## 10. Gain- and loss functions define uncertainty measures

10.1. **Gain- and loss functions.** Although Def. 9.1 of uncertainty measures is abstract, they can be made concrete via "gain functions" [8] or equivalently "loss functions" [4, Eqn. (5)] that encode an attacker's (e.g.) economic interest in the secrets and the cost of obtaining them. We use loss functions here.

**Definition 10.1** (Loss function determines uncertainty measure). A *loss function* $\ell$ is of type $I \to \mathcal{X} \to \mathbb{R}^{\geq}$ for some index set $I$, with the intuitive meaning that $\ell.i.x$ is the cost to the attacker of using "attack strategy" $i$ when the hidden value turns out actually to be $x$. Her expected cost for an attack planned but not yet carried out is then $\mathcal{E}_{\delta}\,(\ell.i)$ if $\delta$ is the distribution in $\mathbb{D}\mathcal{X}$ she believes to be governing $x$ currently.

From such an $\ell$ we define an uncertainty measure

$$U_\ell.\rho \quad := \quad \inf_{i\,:\,I}\ \mathcal{E}_\rho\,(\ell.i)\ . \qquad\qquad (10.1)$$

When $I$ is finite, the inf can be replaced by min.

The inf represents a rational strategy of minimising cost or risk, and a typical attacker will act as follows: she chooses the attack strategy (i.e. he chooses $i$) whose expected cost $\mathcal{E}_\rho\,(\ell.i)$ to her, where $\rho$ is the posterior in $\mathbb{D}\mathcal{X}$ she infers from her observations in $\mathcal{Y}$, will be the least.

**Lemma 10.2** (Well-definedness for Def. 10.1). *For any loss function $\ell\colon I \to \mathcal{X} \to \mathbb{R}^{\geq}$ the function $U_\ell$ in Def. 10.1 is continuous and concave.*

*Proof.* We give here the proof for the finite-$I$ case. (The infinite case is considered in [7, Sec III-B]; it might require further assumptions on $I$.) Let $\ell$ be a loss function and $U_\ell$ be the associated uncertainty measure.

$\underline{U_\ell \text{ is concave}}$: Take $\rho_{1,2}\colon \mathbb{D}\mathcal{X}$ and $p\colon [0,1]$. We have

$$
\begin{array}{lll}
 & U_\ell.(\rho_{1p}+\rho_2) & \\
= & \min_{i:I}\left(\mathcal{E}_{\rho_{1p}+\rho_2}\,\ell.i\right) & \text{"definition } U_\ell\text{"} \\
= & \min_{i:I}\left(\mathcal{E}_{\rho_1}\,\ell.i_p+\mathcal{E}_{\rho_2}\,\ell.i\right) & \text{"}\lambda\delta\cdot\mathcal{E}_\delta\,u\text{ is linear"} \\
\geq & \left(\min_{i:I}\mathcal{E}_{\rho_1}\,\ell.i\right)_p+\left(\min_{i:I}\mathcal{E}_{\rho_2}\,\ell.i\right) & \text{"}(\min f)_p+(\min g)\leq\min(f_p+g)\text{"} \\
= & U_\ell.\rho_{1p}+U_\ell.\rho_2\ . & \text{"definition } U_\ell\text{"}
\end{array}
$$

$\underline{U_l \text{ is continuous}}$: Since $I$ is finite and each function $(\lambda\rho\cdot\mathcal{E}_\rho\,\ell.i)=(\lambda\rho\cdot\sum_{x:\mathcal{X}}\rho_x\times\ell.i.x)$ is continuous, the function $U_\ell$ is also continuous. $\qquad\square$

Remarkably, loss functions are *complete* for uncertainty measures: any uncertainty measure in $\mathbb{U}\mathcal{X}$ can be expressed as $U_\ell$ for some loss function $\ell$ in $I\to\mathcal{X}\to\mathbb{R}^{\geq}$, but possibly requiring $I$ to be infinite [27]. Roughly speaking, this is because of the way concave functions can be expressed as the envelope of their tangential hyperplanes: the coefficients of the hyperplanes' normals are the loss functions. [16]

It is compellingly shown elsewhere how versatile loss (equiv. gain) functions are [8]. Of particular interest is that Lem. 9.2 applies, both in the discrete [4] and the continuous cases [5], even when uncertainties are restricted to those generated by loss functions: the "distinguishing witness" constructed for completeness is in fact a loss function [4].

## 11. A *UM*-transformer example for §4.4

11.1. **Profiling an attacker with a loss function.** In the context of Fig. 4 we imagine an attacker whose livelihood depends on her guessing whether `xs[0]`=`xs[1]` or not, finally. If he guesses incorrectly he loses \$1; if correctly, he breaks even (loses \$0). This is as much a mathematical- as a *social* issue: attacks will be discouraged if they are not worthwhile for the attacker in terms of her own criteria. (See also §4.4 for this social aspect.)

In this example, following §10, we express the attacker's criteria as two strategies "guess same" and "guess different" (thus $I=\{\mathsf{same},\mathsf{diff}\}$) and a loss function $\ell$ therefore defined

$$
\begin{array}{rcllrclc}
 & \ell.\mathsf{same}.(00) & = & 0 & \ell.\mathsf{diff}.(00) & = & 1 & \\
\dagger & \ell.\mathsf{same}.(01) & = & 1 & \ell.\mathsf{diff}.(01) & = & 0 & \ddagger \\
 & \ell.\mathsf{same}.(10) & = & 1 & \ell.\mathsf{diff}.(10) & = & 0 & \\
 & \ell.\mathsf{same}.(11) & = & 0 & \ell.\mathsf{diff}.(11) & = & 1\ , &
\end{array}
$$

based on the informal description just above: for example if `xs=01` but he guesses **same**, the case indicated by $\dagger$, then he loses \$1; but if he guesses **diff**, he breaks even $\ddagger$. Using (10.1) we define our *UM* as $u.\delta=U_\ell.\delta=$

$$
\begin{array}{rccc}
 & \ell.\mathsf{same}.\delta & \min & \ell.\mathsf{diff}.\delta \\
= & \mathcal{E}_\delta\,(\ell.\mathsf{same}) & \min & \mathcal{E}_\delta\,(\ell.\mathsf{diff}) \\
= & (\delta_{00}+\delta_{11}) & \min & (\delta_{01}+\delta_{10})\ .
\end{array}
$$

---

[16]For example, Shannon entropy requires infinite $I$, and the encoding is then related to minimising the Kullback-Leibler divergence.

11.2. **Using *UM*'s and transformers to plan an attack.** We can use our transformer semantics to answer $u$-dependent questions about Fig. 4 over *all* priors: we use the two we chose earlier in §4.4 as examples.

Writing $[\![P]\!]$ for the abstract *HMM* denoted by the two lines of code in Fig. 4, we have for any $\pi$ that

$$\text{wp.}[\![P]\!].u.\pi \quad = \quad \begin{aligned} & \pi_{00} \ \min \ (\pi_{01}+\pi_{10})/2 \\ + \ & \pi_{11} \ \min \ (\pi_{01}+\pi_{10})/2 \ . \end{aligned} \tag{11.1}$$

(See §I below for how this wp.$(\cdot)$ is calculated.)

Now let $\pi^4$ be the prior described by the initial comment in Fig. 4. The attacker's (expected) uncertainty wrt. the *final* hyper $[\![P]\!].\pi^4$ is given by wp.$[\![P]\!].u$ applied to that *initial* (uniform) prior $\pi^4$, that is wp.$[\![P]\!].u.\pi^4 = 1/2$ directly from (11.1). Since $u.\pi^4$ is also $1/2$, he is indifferent wrt. whether he should attack before or after $P$ has been allowed to run.

Now suppose that $\mathtt{xs}[0]=1$ is known initially, thus with prior $\pi$ being $(0,0,1/2,1/2)$ so that $u$ applied initially gives $u.\pi{=}1/2$. But $u$ applied finally would give wp.$[\![P]\!].u.\pi = (0 \min 1/4) + (1/2 \min 1/4) = 1/4 < 1/2$, so that it is better to attack later even though $\mathtt{xs}$ might have been altered by $P$. This scenario confirms that in fact for some priors, the program in Fig. 4 cannot be regarded as secure.

## 12. *HMM*'s and the Dalenius Desideratum

Our abstracting from initial-state correlations allows a semantics for programs' *final* states alone. Sometimes, however, leakage from the *initial* state is important, even if that state is overwritten by the markov part of the *HMM*: what the initial state *was* might reveal information about what some other correlated state still *is*, even if that other state is not mentioned in the program at all. This general concern was raised wrt. statistical databases by Dalenius [10] who argued that it is inescapable; Dwork later gave a proof of this [11]. Here is an (edited) extract from her paper:

> Suppose we have a statistical database that [records] average heights of population subgroups, and suppose further that it is infeasible to learn this information (perhaps for financial reasons) in any other way (say, by conducting a new study). Finally, suppose that one's true height is considered sensitive.
>
> [An adversary having the] auxiliary information "Turing is two inches taller than the average Lithuanian woman" [would, with access to the database, learn] Turing's height. In contrast, anyone without access to the database, knowing only the auxiliary information, learns much less about Turing's height.

With our constructions here, we are able to see the Dalenius effect in programming terms. The program that allows access to Dwork's database of Lithuanian heights (as above) might itself, in isolation, have been analysed for security leaks. But if that program is run in a larger programming context in which there is a (to be kept secret) variable $\mathtt{tHeight}$ containing Turing's height *and* there is program code (external to the database-access code) that establishes a correlation between the two, then running the database-access program reveals information about $\mathtt{tHeight}$ even though that variable is not mentioned anywhere in the database program.

In more austere terms, we would explain the effect as follows. A "classical" sequential program does not affect variables to which it does not refer; for example $\mathtt{x:= E}$ does not affect some other variable $\mathtt{y}$ in any way. But the program $\mathtt{leak\ x}$ (recalling the notation

of Fig. 3) can affect *what we know* about variable y even though the program `leak x` does not refer to y at all.

Consider for example an input distribution $(0,0)@1/2, (1,1)@1/2$ on two variables (x,y). Its y-marginal distribution is uniform on $\{0,1\}$. But the output hyper of that program, projected onto y, is $[0]@1/2, [1]@1/2$, showing that the distribution on y is now a point, no longer uniform: [17] with probability $1/2$ that point will be $[0]$, and with probability $1/2$ that point will be $[1]$. Reviewing the leaks of x tells us which point distribution on y we have, and we see essentially the Dalenius effect between "database" x, "query" `leak x` and "third-party data" y.

This effect is exacerbated when we include state updates, as we have done with our abstract *HMM*'s here. (Updates were not considered originally by Dalenius or by Dwork.) For then the program `leak x; x:= 0` and the program `x:= 0` have the same abstract-*HMM* semantics on state-space (just) x, but different semantics on state-space x,y. [18] The Dalenius effect has become, in programming terms, a failure of compositionality wrt. unreferenced global variables.

We show in this section how to deal with that: in brief, we include both the initial- and the final values of the state in our semantics. The crucial point is that we do not have to do more than that, in particular that we do not have to consider "all possible third-party data y of any type".

We now address the details. Consider a "constant" overwrite-by-x markov $M^{\times}_{x,x'} = 1 \text{ if } x'=\text{x } else \text{ } 0$ for some fixed x: $\mathcal{X}$. Then $[\![C{:}M^{\times}]\!] = [\![{:}M^{\times}]\!]$ for any channel $C$, because $C$ has no effect on our knowedge of the final state. We know already what it is going to be.

We now adjust the semantics so that leakage from the initial state is accounted for, even if it is subsequently overwritten. Let $C{:}\mathcal{X}{\to}\mathcal{Y}$ be a channel and $M{:}\mathcal{X}{\to}\mathcal{X}$ a markov, as usual, and let $\mathcal{Z}$ be fresh. Write $C_{\times\mathcal{Z}}$ in $(\mathcal{X}{\times}\mathcal{Z}){\to}\mathcal{Y}$ for the expanded channel

$$(C_{\times\mathcal{Z}})_{(x,z),y} \quad := \quad C_{x,y} \ ,$$

i.e. that $C$ ignores $z$. Similarly $M_{\times\mathcal{Z}}{:}(\mathcal{X}{\times}\mathcal{Z}){\to}(\mathcal{X}{\times}\mathcal{Z})$ is given by

$$(M_{\times\mathcal{Z}})_{(x,z),(x',z')} \quad := \quad M_{x,x'} \text{ if } z{=}z' \text{ else } 0 \ ,$$

i.e. so that $M$ does not change $x$. Thus these definitions ensure that for any $\pi{:}\mathbb{D}(\mathcal{X}{\times}\mathcal{Z})$ neither $\pi{\triangleright}(C_{\times\mathcal{Z}})$ nor $\pi{\triangleright}(M_{\times\mathcal{Z}})$ depends on the $\mathcal{Z}$ component. Take for example $\mathcal{Z}{:=}\{z_0, z_1\}$ consider $C, M$ as below:

$$
C \ = \ \begin{array}{c} \\ x_0{:} \\ x_1{:} \end{array} \begin{array}{cc} y_0 & y_1 \\ \begin{pmatrix} 1 & 0 \\ 1/4 & 3/4 \end{pmatrix} \end{array}
\qquad
M \ = \ \begin{array}{c} \\ x_0{:} \\ x_1{:} \end{array} \begin{array}{cc} x_0 & x_1 \\ \begin{pmatrix} 1/2 & 1/2 \\ 1/2 & 1/2 \end{pmatrix} \end{array}
$$

$$
C_{\times\mathcal{Z}} \quad = \quad \begin{array}{c} (x_0, z_0){:} \\ (x_0, z_1){:} \\ (x_1, z_0){:} \\ (x_1, z_1){:} \end{array} \begin{array}{cc} y_0 & y_1 \\ \begin{pmatrix} 1 & 0 \\ 1 & 0 \\ 1/4 & 3/4 \\ 1/4 & 3/4 \end{pmatrix} \end{array}
$$

---

[17]Since the output is a hyper, if knowledge of y were unaffected we would have the point hyper on the uniform distribution, that is $[0@1/2, 1@1/2]$.

[18]On state-space x, both programs produce the output hyper $[\![0]\!]$ that denotes "x is certainly 0." On x,y however, the first might reveal something about y while the second cannot.

$$
M_{\times \mathcal{Z}} \quad = \quad
\begin{array}{c}
\\
x_0 z_0: \\
x_0 z_1: \\
x_1 z_0: \\
x_1 z_1:
\end{array}
\begin{array}{cccc}
x_0 z_0 & x_0 z_1 & x_1 z_0 & x_1 z_1
\end{array}
\left(
\begin{array}{cccc}
^1/_2 & 0 & ^1/_2 & 0 \\
0 & ^1/_2 & 0 & ^1/_2 \\
^1/_2 & 0 & ^1/_2 & 0 \\
0 & ^1/_2 & 0 & ^1/_2
\end{array}
\right) .
$$

The definitions above show that in $C_{\times \mathcal{Z}}$ the rows of the original $C$ are each repeated $2 = \#\mathcal{Z}$ times; and the subsequent update by $M_{\times \mathcal{Z}}$ leaves $\mathcal{Z}$ unchanged. Observe that these definitions now account for information flows with respect to initial distributions $\mathbb{D}(\mathcal{X} \times \mathcal{Z})$ where, crucially, the $\mathcal{Z}$ component is merely "carried along". But it captures the Dalenius effect mentioned, as we now explain.

Consider an initial distribution $\pi \colon \mathbb{D}(\mathcal{X} \times \mathcal{Z})$ such that $\pi_{x_i, z_j} = 1$ if and only if $i = j$, i.e. that $z$ is a copy of $x$'s initial value. We see that, even though $\mathcal{Z}$ is not accessed by the program at all, if ever $y_1$ is observed then the $\mathcal{Z}$ component must certainly be $z_1$, and if $y_0$ is observed then it is 4 times more likely to be $z_0$ than $z_1$.

Although $\mathcal{Z}$ is arbitrary, it can be shown that this Dalenius effect on any $\mathcal{Z}$ can be determined by the *HMM* semantics *specifically in the case where $\mathcal{X} = \mathcal{Z}$* as just above. That is, we do not have to consider "all $\mathcal{Z}$'s", which would be impractical. Note the construction of a fully compositional semantics for programs with hidden states is requires further extensive conceptual and technical work which we have developed elsewhere [28].

## 13. Related work

There is great diversity in approaches to information flow in (probabilistic) programs, which we have surveyed in our own earlier work [4–7]. Here we have concentrated on general techniques for semantic constructions, in particular those based on monads, duality and refinement.

Refinement of probabilistic programs appeared in [29] where evaluations were used to construct a powerdomain for probabilistic but possibly non-terminating computations; this was extended to include demonic choice in the discrete case in [22, 24], and was significantly generalised in [30]. Our "uncertainty refinement" that combines information flow with functional properties first appeared for information flow in straight-line programs in [4], was extended to general measure spaces [5] and appeared independently for the specific case of channels [8]. Whereas Jones and Plotkin began with an underlying partial order over which to construct a probability space, our uncertainty-refinement order begins "one level up", using hyper-distributions $\mathbb{D}^2 \mathcal{X}$ to encode an "attack model" that accounts for information flow.

Doberkat defines *stochastic relations* that correspond to forward-semantic functions of type $\mathcal{X} \to \mathbb{D}\mathcal{X}$ for Markov processes: these are what we generalise by going "one level up". The converse of those stochastic relations [31] might improve the presentation of our Def. 2.7, where a hyper is extracted from a channel and a prior, i.e. from a joint distribution.

Dual models for program semantics include [20], then for probabilistic programs [21, 23] in the purely probabilistic case. Subsequently [22] added demonic choice. And [30, 32] study dual models for probability and nondeterminism using a version of Riesz's representation theorem.

In particular, Goubault-Larrecq's approach [32] to combining probability and nondeterminism differs from our earlier work [22]. It uses general denotations for probabilistic programs in which nondeterminism is introduced at the level of measures (by weakening the
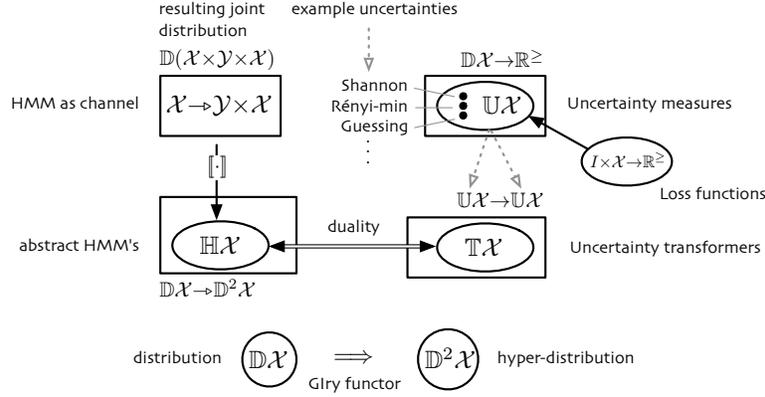
FIGURE 6. Relationship between the semantic spaces.

modularity law) rather than as healthy sets of measures [22, 24, 30]. That leads naturally to a backward semantics of probabilistic demonic programs because nondeterminism is captured within integration. There is thus a strong analogy between our *UM*-transformers and Goubault-Larrecq's "previsions" because both are continuous functionals that act on some set of tests (bounded continuous functions). The main difference is that our *UM*-transformers are specifically tailored to capture security semantics, which is what leads to concavity on our set of uncertainty measures. Notice moreover that Goubault-Larrecq encounters a difficulty similar to our composition of *HMM*'s, that the decomposition $[\![C{:}M]\!]$ (resp. collinearity) is not preserved by Giry composition. Indeed, both difficulties are resolved by working in a larger space, namely, the space of abstract *HMM*'s (resp. not-necessarily-collinear continuous previsions).

In [33] a dual model for Markov processes is used to prove properties about approximations of finite behaviours, and in [34] it is shown how expectation transformers relate to explicit program models described by Markov processes.

Recently Jacobs and Hasuo have explored a general categorical construction of a backward transformer semantics from a forward monadic model of probabilistic computations (discrete, continuous and quantum) [35, 36]. Their construction uses measurability as the underlying feature of "predicates", while the stronger condition of continuity is crucial for our uncertainty measures. It would be interesting to see whether an instantiation of that categorical derivation can provide more structure for what we have done.

The concave functions advocated here for analysing information-flow properties have appeared in [5, 8] and have been identified in [37] as an ingredient in privacy analysis.

## 14. CONCLUSIONS AND PROSPECTS

Our principal objective was to provide an abstract setting for *HMM*'s based on well understood principles of semantic spaces. We did that using Giry's general monadic framework applied at the level of $\mathbb{D}\mathcal{X}$ (rather than $\mathcal{X}$); the resulting structures include a refinement order which is sensitive to both functional *and* information-flow properties, and they lead to a dual, transformer space supported by theorems demonstrating the duality. Fig. 6 summarises the results:

- Top-left Fig. 6 shows the three-way joint distribution $\mathbb{D}(\mathcal{X} \times \mathcal{Y} \times \mathcal{X})$ produced by an *HMM* applied to a prior of type $\mathbb{D}\mathcal{X}$: recall §3.1.
- At bottom-left we formulate *abstract HMM's* over a state $\mathcal{X}$, a type $\mathbb{H}\mathcal{X}$, as a monadic model for *HMM*'s over $\mathcal{X}$, and give their characteristic properties: recall §3.2,6,8.
- At top-right we have *uncertainty measures* $\mathbb{U}\mathcal{X}$ as a generalisation of diverse entropies (top centre), and we gave their characteristic properties: recall §9.1.
- We showed (centre right) that uncertainty measures **have a complete representation** as *loss functions*; recall §10.
- We gave a dual, uncertainty-transformer semantics $\mathbb{T}\mathcal{X}$ of $\mathbb{H}\mathcal{X}$, stating its characteristic properties (bottom right) and proved that they enable the duality with $\mathbb{H}\mathcal{X}$ (centre): recall Thm. 9.8 in §9.2.
- We showed how all of that is an instance of the general Giry monad as a computation, of which (finite) *HMM*'s use a discrete portion (bottom centre): recall §7.
- We explained how the "Dalenius effect" is manifested as a compositional issue in this framework, and suggested how it can be treated: recall §12.
- We stated and proved Thm. 9.8, which we believe is a significant new result, in particular its assumptions and proof.

More abstractly (recall §1.2), we aimed to profit by joining two ideas: the established use of *HMM*'s as descriptions of probabilistic mechanisms having hidden state, and the established use of monads for modelling computations. Our novel use of $\mathbb{D}\mathcal{X}$ in the monad, rather than the state $\mathcal{X}$ itself, is the principal innovation that allowed this; and the synthesised hyper-distribution space that results leads to other advantages (the two †'s below).

An immediate benefit accrues because, in monad-enabled programming languages, probabilistic-programming packages can be built very quickly and e.g. [38] is just one of many examples. Indeed the translation into real programs is almost elementary because of the powerful and general structures available: the Haskell prototype independently verifies the examples in Figs. 2–5. (See §5 for an overview.)

More importantly, any monad brings with it both general equational properties and specific properties applying to the monad in question (such as those in [2]). These conceptual tools allow reasoning about the structures modelled (*HMM*'s in this case) in ways that would be obscured by their more direct operational representation (e.g. as matrices).

- The other advantages of hypers are several: one is that they abstract from differences between entropies in a way that allows all of the entropies to be used uniformly. For example, a hyper contains all the information necessary to calculate the information leakage of a particular program fragment (typically, in the security literature, a pure channel §3.4), as shown in [6], and furthermore the Kantorovich-metric structure of $\mathbb{D}\mathcal{X}$ we used earlier for channels [7] now carries over to *HMM*'s.
- Another advantage of hypers is that their partial-order enables semantics for "looping *HMM*'s" in the standard way (least fixed-point) for computer science, rather than a direct ad-hoc definition based on matrices. Indeed a typical use of *HMM*'s is to run a single *HMM*-step (§3.1) repeatedly and then to make statistical deductions about its hidden features: sophisticated mathematical tools are available for this special case [12]. Via abstract *HMM*'s we can however, in principle, handle complex, heterogeneous systems beyond (what amounts to, in the special case just above) a single loop containing just a single statement.

Our more concrete aim (again §1.2) was to allow source-level reasoning about probabilistic programs with hidden state. Historically *at the source level* this works best with backwards reasoning based on predicates (or similar) that can be embedded between program statements rather than forwards reasoning which, here, would be calculations using $\mathbb{D}\mathcal{X} \to \mathbb{D}^2\mathcal{X}$ directly.

Here our "predicates" are *UM*'s, which in this paper however are mathematical objects unsuitable for embedding directly in program texts (see §I, last paragraph) As remarked in §10.1, however, any *UM* can be expressed as $U_l$ for some loss-function $l$ which function –crucially– is indeed an expression based on program variables [27]. The added complexity introduced by the hidden state is that the program-logic based on that observation must represent the index-set ($I$) of the loss function; that would most likely be done by adding a special-purpose quantifier (since the loss-function index must be a bound variable within the assertion, not appearing in the program proper).

Exploiting this opportunity for a source-level quantitative logic of probabilistic hidden state is planned for future work.

## REFERENCES

[1] E. Moggi, "Computational lambda-calculus and monads," in *Proc. 4th IEEE Symp. LiCS*, 1989, pp. 14–23.

[2] M. Giry, "A categorical approach to probability theory," in *Categorical Aspects of Topology and Analysis*, ser. Lecture Notes in Mathematics. Springer, 1981, vol. 915, pp. 68–85.

[3] A. McIver, L. Meinicke, and C. Morgan, "Hidden-Markov program algebra with iteration," *Mathematical Structures in Computer Science*, 2014.

[4] ——, "Compositional closure for Bayes risk in probabilistic noninterference," in *Proc. 37th Int. Colloq. ICALP 2010, Part II*, 2010, pp. 223–235.

[5] ——, "A Kantorovich-monadic powerdomain for information hiding, with probability and nondeterminism," in *Proc. 27th Symp. LiCS*, 2012, pp. 460–70.

[6] A. McIver, C. Morgan, G. Smith, B. Espinoza, and L. Meinicke, "Abstract channels and their robust information-leakage ordering," in *Proc. 3rd Conf. PoST (ETAPS)*, ser. Lecture Notes in Computer Science, M. Abadi and S. Kremer, Eds., vol. 8414. Springer, 2014, pp. 83–102.

[7] M. S. Alvim, K. Chatzikokolakis, A. McIver, C. Morgan, C. Palamidessi, and G. Smith, "Additive and multiplicative notions of leakage, and their capacities," in *Proc 27th IEEE Symp. CSF*. IEEE, 2014, pp. 308–322.

[8] M. S. Alvim, K. Chatzikokolakis, C. Palamidessi, and G. Smith, "Measuring information leakage using generalized gain functions," in *Proc. 25th IEEE Symp. CSF*, Jun. 2012, pp. 265–79.

[9] G. Smith, "On the foundations of quantitative information flow," in *Proc. 12th Conf. FoSSaCS (ETAPS)*, ser. Lecture Notes in Computer Science, L. de Alfaro, Ed., vol. 5504, 2009, pp. 288–302.

[10] T. Dalenius, "Towards a methodology for statistical disclosure control," *Statistik Tidskrift*, vol. 15, pp. 429–44, 1977.

[11] C. Dwork, "Differential privacy," in *Proc. 33rd Int. Colloq. ICALP*, 2006, pp. 1–12.

[12] D. Jurafsky and J. Martin, *Speech and Language Processing.* Prentice Hall International, 2000.

[13] T. M. Cover and J. A. Thomas, *Elements of Information Theory*, 2nd ed. John Wiley & Sons, Inc., 2006.

[14] D. Fremlin, *Measure Theory.* Torres Fremlin, 2000.

[15] J. Gibbons, A. McIver, C. Morgan, and T. Schrijvers, "Quantitative information flow with monads in haskell," 2018, submitted for review.

[16] J. Landauer and T. Redmond, "A lattice of information," in *Proc. 6th IEEE CSFW'93*, Jun. 1993, pp. 65–70.

[17] F. van Breugel, "The metric monad for probabilistic nondeterminism," 2005, `www.cse.yorku.ca/∼franck/research/drafts/monad.pdf`.

[18] K. R. Parthasarathy, *Probability Measures on Metric Spaces.* AMS Chelsea Publishing, 1967.

[19] A. L. Gibbs and F. E. Su, "On choosing and bounding probability metrics," *International Statistical Review*, vol. 70, no. 3, pp. 419–435, 2002. [Online]. Available: http://dx.doi.org/10.1111/j.1751-5823.2002.tb00178.x

[20] E. Dijkstra, *A Discipline of Programming.* Prentice-Hall, 1976.

[21] D. Kozen, "A probabilistic PDL," in *Proc. 15th ACM Symp. Theory of Computing.* ACM, 1983, pp. 291–7.

[22] C. Morgan, A. McIver, and K. Seidel, "Probabilistic predicate transformers," *ACM Trans Prog Lang Sys*, vol. 18, no. 3, pp. 325–53, 1996.

[23] C. Jones, "Probabilistic nondeterminism," Edinburgh University, Monograph ECS-LFCS-90-105, 1990, (Ph.D. Thesis).

[24] A. McIver and C. Morgan, *Abstraction, Refinement and Proof for Probabilistic Systems*, ser. Tech Mono Comp Sci. Springer, 2005.

[25] M. H. Stone, "The generalized Weierstrass approximation theorem," *Math Magazine*, vol. 21, no. 4, pp. 167–184, March 1948.

[26] R. Ranga Rao, "Relations between weak and uniform convergence of measures with applications," *Annals of Mathematical Statistics*, vol. 33, no. 2, pp. 659–680, January 1962.

[27] K. Chatzikokolakis, Private communications, 2014.

[28] A. McIver, C. C. Morgan, and T. M. Rabehaja, "Algebra for quantitative information flow," in *Relational and Algebraic Methods in Computer Science - 16th International Conference, RAMiCS 2017, Lyon, France, May 15-18, 2017, Proceedings*, 2017, pp. 3–23. [Online]. Available: https://doi.org/10.1007/978-3-319-57418-9_1

[29] C. Jones and G. Plotkin, "A probabilistic powerdomain of evaluations," in *Proc. 4th IEEE Symp. LiCS*, 1989, pp. 186–95.

[30] R. Tix, K. Keimel, and G. Plotkin, "Semantic domains for combining probability and non-determinism," *Electron. Notes Theor. Comput. Sci.*, vol. 222, pp. 3–99, 2009.

[31] E. Doberkat, "The converse of a stochastic relation," in *Proc. 6th Conf. FoSSaCS (ETAPS)*, ser. LNCS, A. Gordon, Ed., vol. 2620. Springer-Verlag, 2003, pp. 233–49.

[32] J. Goubault-Larrecq, "Continuous previsions," in *Proc. 16th EACSL*, ser. Lecture Notes in Computer Science, vol. 4646. Springer, 2007, pp. 542–57.

[33] P. Chaput, V. Danos, P. Panangaden, and G. D. Plotkin, "Approximating Markov processes by averaging," *J. ACM*, vol. 61, no. 1, 2014.

[34] F. Gretz, J. Katoen, and A. McIver, "Operational versus weakest pre-expectation semantics for the probabilistic guarded command language," *Perform. Eval.*, vol. 73, pp. 110–132, 2014.

[35] B. Jacobs, "Measurable spaces and their effect logic," in *Proc. 28th LiCS*, 2013, pp. 83–92.

[36] I. Hasuo, "Generic weakest precondition semantics from monads enriched with order," in *Proc. CMCS*, ser. LNCS, M. Bonsangue, Ed., vol. 8446. Springer, 2014, pp. 10–32.

[37] D. Kifer and B.-R. Lin, "Towards an axiomatization of statistical privacy and utility," 2010, Penn State Technical report: CSE-10-002.

[38] M. Erwig and S. Kollmansberger, "Probabilistic functional programming in Haskell," *Journal of Functional Programming*, vol. 16, pp. 21–34, 2006.

[39] D. Blackwell, "The comparison of experiments," in *Proc. 2nd Berkely Symp. Mathematical Statistics and Probability.* Univ. Califormia Press, 1951, pp. 93–102.

[40] M. Bačák and J. M. Browein, "On difference convexity of locally Lipschitz functions," *Optimization: A Journal of Math Prog and Oper Research*, vol. 60, no. 8-9, pp. 961–978, 2011.

[41] L. H. Loomis and S. Sternberg, *Advanced Calculus.* Jones and Bartlett Publishers, 1990.

## Appendix A. Summary of notations

These entries list in first-use order the points at which notation is introduced during the exposition: a detailed explanation of each is given there.

| | | |
|---|---|---|
| $-^{\dagger}$ | Kleisli extension | p.2 |
| $f.x$ vs. $f(x)$ | Function application is ".", i.e. a dot. | p.3 |
| $\mathcal{X} \rightarrow \mathcal{Y}$ | Type of a matrix. | p.3 |
| $C_{x,y}, C_{-,y}$ etc. | Elements of vectors and matrixes by index; whole rows/-columns. | p.3 |
| $\vec{\mathcal{X}}$ | Type of vector. | p.3 |
| $(\cdot)$ | Matrix multiplication: vectors automatically taken as row- or column- for conformity. | p.3 |
| $(:)$ vs. $(\in)$ | Declaration vs. property. | p.3 |
| u.c. Roman letter | Matrices: $C$ for channels; $M$ for transformers; $H$ for $HMM$'s | p.4 |
| $\mathcal{X}$ | Finite set of states. | p.4 |
| $\mathcal{Y}$ | Finite set of observations. | p.4 |
| l.c. Greek letter | Vectors, usually distributions over $\mathcal{X}$: $\pi$ for priors; $\rho$ for posteriors; $\delta$ for others. | p.4 |
| $\Sigma()$ | Weight (sum of elements) of vector or matrix. | p.4 |
| $\pi \triangleright C$ | Channel applied to prior. | p.4 |
| $\lfloor - \rfloor$ | Normalisation of distribution. | p.4 |
| *similar* | wrt. columns of joint matrix. | p.5 |
| $y_{1,2}$ vs. $y_1, y_2$ | Former abbreviates latter. | p.5 |
| $\mathbb{D}$ | Discrete-distribution type constructor, a functor. | p.5 |
| $\mathbb{D}^2$ | Distribution-of-distributions. | p.5 |
| hyper | Abbreviation of "hyper-distribution". | p.5 |
| inner | Element of a hyper's base type. | p.5 |
| outer | Distribution of a hyper over its inners. | p.5 |
| $[\![\cdot]\!]$ | Semantic function for $HMM$'s. | p.6 |
| $\underline{\mathbb{D}}$ | Sub-distribution. | p.6 |
| $\underline{\mathbb{D}}^2$ | Sub-hyper. | p.6 |
| $x_p \oplus x'$ | The two-point distribution "$x$ with probability $p$ and $x'$ with probability $1-p$". | p.6 |
| $[\cdot]$ | Point distribution. | p.6 |
| $[-]$ | Sub-point distribution. | p.7 |
| $\lceil \delta \rceil$ | The support of a distribution. | p.7 |
| avg | Average (of hyper); multiply in monad. | p.7 |
| $\Delta$ | Upper-case Greek for hypers. | p.7 |
| channel | The emission part of an $HMM$-step. | p.8 |
| markov | The transition part of an $HMM$-step. | p.8 |
| $(C;M), (C;), (;M)$ | One-step $HMM$ defined by channel and markov. | p.8 |
| nc | Channel that releases no information. | p.9 |
| id | Identity (Markov) transform. | p.10 |
| @ | Notation for specific hyper-distributions | p.10 |

## Appendix B. Characterisation of pure channels and pure markovs                 [§4.1]

B.1. **Pure abstract markovs.** Since a pure markov reveals nothing, a pure abstract markov $h\colon \mathbb{H}\mathcal{X}$ should produce only point hypers, i.e. have for all $\pi\colon \mathbb{D}\mathcal{X}$ that $h.\pi = [\rho]$ for some $\rho$ (depending on $\pi$).

From that we can deduce that for any pure abstract markov $h$ the effect of $\mathsf{avg} \circ h$ (on some $\pi$) is matrix multiplication by some $M$ (independent of $\pi$). That is, for any $0 \le p \le 1$ we have

$$(\mathsf{avg} \circ h).(\pi_{1p}+\pi_2) \quad = \quad (\mathsf{avg} \circ h).\pi_{1\,p}+ (\mathsf{avg} \circ h).\pi_2 \ , \tag{B.1}$$

which property characterises matrix multiplication. This is because $h.(\pi_{1p}+\pi_2)=[\rho]$ and $h.\pi_{1,2}=[\rho_{1,2}]$ resp. for some $\rho, \rho_{1,2}$, together with Lem. 8.3, gives

$$[\rho_1]_p+[\rho_2] \quad \sqsubseteq \quad [\rho] \ ,$$

and the only way that can hold is if $\rho = \rho_{1p}+\rho_2$, which is precisely the claim made at (B.1) just above.

B.2. **Pure abstract channels.** A pure channel is one that releases information about the distribution on $\mathcal{X}$ but does not change it: one can think of the transformation part as the identity matrix. Thus (B.1) above suggests that we should have that $\mathsf{avg} \circ h$ is the identity for a pure channel, i.e. that $\mathsf{avg}.(h.\pi){=}\pi$. This is necessary, but turns out not to be sufficient: we explore a fuller characterisation of channels later (§J).

APPENDIX C. EQUIVALENT PRESENTATIONS OF REFINEMENT: LEM. 6.2               [§6]

Lem. 6.2 concerned two definitions of uncertainty refinement, showing them to be equivalent: one was formulated for joint distributions (defined at (6.1) within the lemma), suitable for discrete reasoning; and the other was formulated for hypers (Def. 6.1), suitable for extension to more general reasoning (e.g. proper measures). We sketch the proof of that equivalence in §D immediately below.

   In this section however we present an example, two hypers $\Delta_{S,I}$ shown to satisfy $\Delta_S{\sqsubseteq}\Delta_I$ in both presentations (Def. 6.1, Lem. 6.2), with an explanation of how to move from one presentation to the other.

   As in (3.1) of §3.4, we use the following notation for discrete distributions where specific values in the support are named: we write

$$
\begin{array}{l}
x_1 \,@\, p_1 \\
x_2 \,@\, p_2 \\
etc...
\end{array}
\tag{C.1}
$$

If these are laid out horizontally, we enclose them in double set-brackets $\{\{\cdots\}\}$ separated by commas: thus $\{\{H@2/3, T@1/3\}\}$ describes a coin twice as likely to give heads as tails. If the double brackets are used without probabilities (and thus also without @'s) then the intended distribution is uniform, so that $\{\{H, T\}\}$ describes a fair coin; a convenient special case of that is $\{\{H\}\}$ for the point distribution on $H$, the coin that gives heads every time. [19]

   Let $\mathcal{X}$ be the set $\{H, T\}$ of coin-flip results. We choose our two hypers as follows, presenting them as at (C.1):

$$
\Delta_S \;\; = \;\; \left[ \begin{array}{ll}
H_{2/3}{\oplus}T & @\,1/2 \\
H_{1/3}{\oplus}T & @\,1/2
\end{array} \right]
$$

$$
\Delta_I \;\; = \;\; \left[ \begin{array}{ll}
H_{2/3}{\oplus}T & @\,1/3 \\
H_{1/2}{\oplus}T & @\,1/3 \\
H_{1/3}{\oplus}T & @\,1/3 \; .
\end{array} \right]
$$

   The first hyper $\Delta_S$ represents choosing fairly between two biased coins and having the chosen one secretly flipped: we know which coin was flipped, but we are not allowed to see the outcome of the flip. In $\Delta_I$ however we choose fairly between *three* coins: the two biased coins from before, and a fair one. Again the chosen one is secretly flipped; again we are not allowed to see the outcome.

   We argue that in any reasonable measure of secrecy, it should in the second case $\Delta_I$ be harder to guess which of $H, T$ resulted from the flip than in the first case $\Delta_S$. And it is

---

[19] In the semantic space we write $[x]$ for that: here we are syntactic.

precisely that non-specific "in *any* reasonable measure" that uncertainty refinement $\Delta_S \sqsubseteq \Delta_I$ attempts to capture. [20]

In this case, and informally speaking, $\Delta_I$ is more secure than $\Delta_S$ because there is now a third possible case that acts as a linear combination of the existing two. That is, some of the separation between the inners $H_{2/3} \oplus T$ and $H_{1/3} \oplus T$ in the support of $\Delta_S$ has been merged together to become a single inner $H_{1/2} \oplus T$ in the support of $\Delta_I$ — and what makes the observer more uncertain is that he doesn't know how to pull that single inner apart again.

Two (reduced) joint matrices $J_{S,I}$ that give $\Delta_{S,I}$ resp. are

$$J_S \quad = \quad \begin{matrix} H: \\ T: \end{matrix} \quad \begin{matrix} a & b \\ \begin{pmatrix} 1/3 & 1/6 \\ 1/6 & 1/3 \end{pmatrix} \end{matrix}$$

$$J_I \quad = \quad \begin{matrix} H: \\ T: \end{matrix} \quad \begin{matrix} c & d & e \\ \begin{pmatrix} 2/9 & 1/6 & 1/6 \\ 1/9 & 1/6 & 2/9 \end{pmatrix} \end{matrix}$$

where the observation spaces are $\mathcal{Y}_S = \{a, b\}$ and $\mathcal{Y}_I = \{c, d, e\}$ respectively. (the column names are arbitrary.) Now the refinement matrix that establishes (according to Lem. 6.2) that $\Delta_S \sqsubseteq \Delta_I$ is $R: \mathcal{Y}_S \rightarrow \mathcal{Y}_I$ given by

$$R \quad = \quad \begin{matrix} a: \\ b: \end{matrix} \quad \begin{matrix} c & d & e \\ \begin{pmatrix} 2/3 & 1/3 & 0 \\ 0 & 1/3 & 2/3 \end{pmatrix} \end{matrix}$$

which, read columnwise, says in its column $c$ that to make Column $c$ of $J_I$ you take $2/3$ of Column $a$ of $J_S$ and none of Column $b$ of $J_S$. The middle column $d$ of $R$ is where the actual refinement lies, that Column $d$ of $J_I$ is made by adding $1/6$ of each of Columns $a, b$ of $J_S$ together. This is where $J_I$ (equiv. $\Delta_I$) reveals less than $J_S$ (equiv. $\Delta_S$) does about the distribution on $\mathcal{X} = \{H, T\}$. And, as the lemma suggests, we indeed have $J_S \cdot R = J_I$.

The alternative, more abstract presentation of this is in terms of Def. 6.1, i.e. where the $\underline{\Delta}$ we are looking for, that establishes $\Delta_S \sqsubseteq \Delta_I$ at the hyper-level directly, can be given as (the denotation of) a joint distribution $J: \mathbb{D}\mathcal{X} \rightarrow \mathcal{Y}_I$ itself: we will have $\underline{\Delta} := [\![J]\!]$ which, because $J$'s source type is $\mathbb{D}\mathcal{X}$, will have type $\mathbb{D}^2(\mathbb{D}\mathcal{X}) = \mathbb{D}^3\mathcal{X}$ as we expect from $[\![\cdot]\!]$. The rows of $J$ will be labelled by the support of $\Delta_S$, i.e. it will have only two rows so that we have

$$J \quad = \quad \begin{matrix} H_{2/3} \oplus T: \\ H_{1/3} \oplus T: \end{matrix} \quad \begin{matrix} c & d & e \\ \begin{pmatrix} 1/3 & 1/6 & 0 \\ 0 & 1/6 & 1/3 \end{pmatrix} \end{matrix}. \tag{C.2}$$

If on the other hand we were to write $\underline{\Delta} = [\![J]\!]$ as a hyper directly (performing the various normalisations etc.) we would have

$$\underline{\Delta} \quad = \quad \begin{bmatrix} [\, H_{2/3} \oplus T \,] & @\, 1/3 \\ (H_{2/3} \oplus T)_{1/2} \oplus (H_{1/3} \oplus T) & @\, 1/3 \\ [\, H_{1/3} \oplus T \,] & @\, 1/3 \end{bmatrix},$$

with each inner here corresponding to a row of (C.2).

Now $\mathsf{avg}.\underline{\Delta}$ is given by the calculation

---

[20] Furthermore, the powerful "Coriaceous" completeness property (Lem. 9.2) shows the dual result: if some $\Delta_S, \Delta_I$ are *not* in the refinement relation, that is $\Delta_S \not\sqsubseteq \Delta_I$, then there is *guaranteed* to be a uncertainty measure wrt. to which $\Delta_I$ is *not* more secure than $\Delta_S$.

$$\begin{array}{ll} & [\ H_{2/3}\oplus T\ ]\times 1/3 \\ + & (H_{2/3}\oplus T)_{1/2}\oplus(H_{1/3}\oplus T)\times 1/3 \\ + & [\ H_{1/3}\oplus T\ ]\times 1/3 \end{array}$$

$$\begin{array}{ll} = & H_{2/3}\oplus T)_{1/2}\oplus(H_{1/3}\oplus T) \\ = & \Delta_S\ . \end{array}$$

This can also be seen (indeed is easier to see) if we simply take the left-marginal of $J$, for which you add the columns together: you get

$$\begin{array}{c} c+d+e \\ = 1 \\ \begin{array}{ll} H_{2/3}\oplus T: & \begin{pmatrix} 1/2 \\ 1/2 \end{pmatrix} \end{array} \end{array}\quad,$$
$$\begin{array}{ll} H_{1/3}\oplus T: \end{array}$$

which is again $\Delta_S$.

For the other direction we obtain $(\mathbb{D}\mathsf{avg}).\underline{\Delta}$ by $\mathsf{avg}$'ing each inner of $\underline{\Delta}$ while preserving the (outer) probabilities. [21] That gives

$$(\mathbb{D}\mathsf{avg}).\underline{\Delta}\quad=\quad \left[\begin{array}{ll} H_{2/3}\oplus T & @\ 1/3 \\ H_{1/2}\oplus T & @\ 1/3 \\ H_{1/3}\oplus T & @\ 1/3 \end{array}\right]\ ,$$

because

$$\begin{array}{rll} \mathsf{avg}.[\ H_{2/3}\oplus T\ ] & = & H_{2/3}\oplus T \\ \mathsf{avg}.(H_{2/3}\oplus T)_{1/2}\oplus(H_{1/3}\oplus T) & = & H_{1/2}\oplus T \\ \mathsf{avg}.[\ H_{1/3}\oplus T\ ] & = & H_{1/3}\oplus T \end{array}$$

And so that the remaining question is "How do we get such a $\underline{\Delta}$ from a given $R$?"

Remember that the support of $\Delta_S$ is $\{H_{2/3}\oplus T, H_{1/3}\oplus T\}$. Make a distribution $\pi_S$ by mapping those (inner) distributions of $\Delta_S$ onto the labels in $\mathcal{Y}_S$ associated uniquely with them in $J_S$. (The association is unique because $J_S$ is reduced.) That gives us that $\pi_S$ is of type $\mathbb{D}\mathcal{Y}_S$ and has value $a_{1/2}\oplus b$.

Now form the joint-distribution matrix $\pi_S\triangleright R$, i.e.

$$\begin{array}{cc} & \begin{array}{ccc} c & d & e \end{array} \qquad\qquad \begin{array}{ccc} c & d & e \end{array} \\ \begin{array}{l} a: \begin{pmatrix} 1/2 \\ 1/2 \end{pmatrix} \end{array} \begin{array}{l} b: \end{array} \triangleright \begin{pmatrix} 2/3 & 2/3 & 0 \\ 0 & 1/3 & 2/3 \end{pmatrix} = \begin{array}{l} a: \\ b: \end{array} \begin{pmatrix} 1/3 & 1/6 & 0 \\ 0 & 1/6 & 1/3 \end{pmatrix} \end{array}$$

which (like $R$ itself) is of type $\mathcal{Y}_S{\rightarrow}\mathcal{Y}_I$. (But note that $R$ is a channel matrix, whereas $\pi_S\triangleright R$ is a joint-distribution matrix.)

Now use the relabelling in the reverse direction on the rows of the joint distribution above (as "new row-labels" at right above) to get a matrix with the same contents but now of type $\mathbb{D}\mathcal{X}{\rightarrow}\mathcal{Y}_I$. It is

$$\begin{array}{cc} & \begin{array}{ccc} c & d & e \end{array} \\ \begin{array}{l} H_{2/3}\oplus T: \\ H_{1/3}\oplus T: \end{array} & \begin{pmatrix} 1/3 & 1/6 & 0 \\ 0 & 1/6 & 1/3 \end{pmatrix} \end{array}$$

which is exactly the $J$ we had at (C.2) above, and as above we get $\underline{\Delta}$ via $\underline{\Delta}=[\![J]\!]$.

---

[21]Recall that the inners of $\underline{\Delta}$ are themselves hypers, which is why they can be $\mathsf{avg}$'d.

Thus in this example we have illustrated how one might move between the two equivalent definitions of refinement. Each one has a witness: in the hyper-formulation it is the distribution on hypers $\underline{\Delta}$; and in the matrix formulation is is a post-processing "refinement matrix" $R$. The sketch proof (§D) shows how to obtain each from the other in general.

## Appendix D. Monadic vs. matrix presentations of refinement       [§6]

In §C we gave an example of the two equivalent presentations of refinement; here we give a proof (sketch) that it can always be done.

*Lemma 6.2: Refinement of joint-distribution matrices*    Let $J_S\colon \mathcal{X}'{\twoheadrightarrow}\mathcal{Y}_S$ and $J_I\colon \mathcal{X}'{\twoheadrightarrow}\mathcal{Y}_I$ be joint-distribution matrices, both of them *reduced* in the sense of Def. 2.7, such that $[\![J_{S,I}]\!]=\Delta_{S,I}$ resp. In this section only we use $\mathcal{X}'$ as a reminder that the *input* side of these $J$'s, their row-indices, is actually the *output* side of the *HMM*'s from which they are derived, i.e. that $J_{x',y} = \sum_x H_{x,y,x'}$ as in Def. 3.2.

We prove the equivalence

$$\Delta_S \sqsubseteq \Delta_I \quad \text{iff} \quad J_S \cdot R = J_I \quad \text{for some } R$$

where $R$ is a stochastic *refinement matrix* of type $\mathcal{Y}_S{\twoheadrightarrow}\mathcal{Y}_I$ (i.e. such that $\sum R_{y,-} = 1$ for each $y\colon \mathcal{Y}_S$).

*Proof.* First we note that for any reduced joint distribution matrix $J\colon \mathcal{Z}{\twoheadrightarrow}\mathcal{Z}'$ there is a one-one correspondence between $J$'s column labels, i.e. elements of $\mathcal{Z}'$, and the support of the hyper $\Delta=[\![J]\!]$ that $J$ defines: it is the function $j\colon \mathcal{Z}' \xleftrightarrow{\text{1-1}} \lceil\Delta\rceil$ from Def. 2.7, injective into $\mathbb{D}\mathcal{Z}$ because $J$ is reduced. We write ($\xleftrightarrow{\text{1-1}}$) to emphasise our one-one use of it below.

$R$ **makes** $\underline{\Delta}$:   We show first that for $J_{S,I}, \Delta_{S,I}$ and $R$ as above we can construct a suitable $\underline{\Delta}$. Let the relabelling associated with $J_S$ be $j_S\colon Y_S \xleftrightarrow{\text{1-1}} \lceil\Delta_S\rceil$. Relabel $\Delta_S$ so that it is a distribution of support $\mathbb{D}Y_S$, so that we can use Def. 2.7 to define $\underline{\Delta}:=[\![\Delta_S{\triangleright}R]\!]$, noting that the types of (relabelled) $\Delta_S\in\mathbb{D}\mathcal{Y}_S$ and of $R\in\mathcal{Y}_S{\twoheadrightarrow}\mathcal{Y}_I$ are precisely what Def. 2.7 requires to produce a result of type $\mathbb{D}^2\mathcal{Y}_S$. Now relabel this (back again) to make an element of $\mathbb{D}^2\lceil\Delta_S\rceil$, that is of $\mathbb{D}^3\mathcal{X}$ because $\lceil\Delta_S\rceil{\subseteq}\mathbb{D}\mathcal{X}$.

We have $\Delta_S=\mathsf{avg}.\underline{\Delta}$ immediately, from the remark following Def. 2.8.

For $(\mathbb{D}\mathsf{avg}).\underline{\Delta}=\Delta_I$ we first calculate

$$
\begin{aligned}
& (\mathbb{D}\mathsf{avg}).[\![\Delta_S{\triangleright}R]\!] \\
={} & [\![M{\cdot}(\Delta_S{\triangleright}R)]\!]\,, \qquad\qquad\qquad \text{``Set } \mathcal{D}:=\lceil\Delta_S\rceil \text{ in Lem. D.1 below''} \\
& \text{where } M_{x,\rho}:=\rho.x \text{ for } x\colon \mathcal{X} \text{ and } \rho\colon \lceil\Delta_S\rceil.
\end{aligned}
$$

Now for arbitrary $x\colon \mathcal{X}$ and $y_I\colon \mathcal{Y}_I$ we continue

$$
\begin{aligned}
& (M{\cdot}(\Delta_S{\triangleright}R))_{x,y_I} \\
={} & \sum_{\rho\colon \mathcal{D}} M_{x,\rho}\,(\Delta_S)_\rho\,R_{\rho,y_I} \\
={} & \sum_{\rho\colon \lceil\Delta_S\rceil} \rho.x\,(\Delta_S)_\rho\,R_{\rho,y_I} & \text{``Defn. } M;\ \mathcal{D}{=}\lceil\Delta_S\rceil\text{''} \\
={} & \sum_{y_S\colon \mathcal{Y}_S} (J_S)_{x,y_S}\,R_{\rho,y_I} & \text{``}\lceil\Delta_S\rceil{=}\mathcal{Y}_S;\ \Delta_S{=}[\![J_S]\!]\text{''} \\
={} & J_I\,, & \text{``}J_I{=}J_S{\cdot}R\text{''}
\end{aligned}
$$

whence $(\mathbb{D}\mathsf{avg}).\underline{\Delta} = [\![J_I]\!] = \Delta_I$ as required.

$\underline{\Delta}$ **makes** $R$:   To show that from $\underline{\Delta}$ we can construct a suitable $R$, we do similar calculations to the above, but in the reverse direction.                                  □

**Lemma D.1** (Technical lemma). *Let $\mathcal{D}{\subseteq}\mathbb{D}\mathcal{X}$ be some finite set of distributions on $\mathcal{X}$, and let $J{:}\mathcal{D}{\rightarrow}\mathcal{Y}$ be a joint-distribution matrix between (those) distributions on $\mathcal{X}$ and some observation space $\mathcal{Y}$. Then $\mathbb{D}\mathsf{avg}.[\![J]\!] = [\![M{\cdot}J]\!]$, where $M{:}\mathcal{X}{\rightarrow}\mathcal{D}$ is defined $M_{x,\delta}{:=}\delta.x$ for $x{:}\mathcal{X}$ and $\delta{:}\mathcal{D}$.*

*Proof.* Let us match inners (and associated weight) of $[\![M{\cdot}J]\!]$ with that of $\mathbb{D}\mathsf{avg}.[\![J]\!]$. These are finitely supported distributions so the following sums are all finite.

Let $y{:}\mathcal{Y}$. On the one hand, the $y$-inner $\delta$ of $[\![M{\cdot}J]\!]$ satisfies, for every $x{:}\mathcal{X}$,

$$\delta.x = \frac{(M \cdot J)_{x,y}}{\sum_x (M \cdot J)_{x,y}} = \frac{\sum_\rho M_{x,\rho} J_{\rho,y}}{\sum_x \sum_\rho M_{x,\rho} J_{\rho,y}} = \frac{\sum_\rho \rho.x J_{\rho,y}}{\sum_\rho J_{\rho,y}}$$

and this inner has weight $\sum_\rho J_{\rho,y}$.

On the other hand, the $y$-inner $\Delta$ of $[\![J]\!]$ satisfies, for every $\rho$,

$$\Delta.\rho = \frac{J_{\rho,y}}{\sum_\rho J_{\rho,y}} \ .$$

This inner has weight $\sum_\rho J_{\rho,y}$. Applying the $\mathsf{avg}$, we get

$$\mathsf{avg}.\Delta.x = \sum_\rho \rho.x \Delta.\rho = \frac{\sum_\rho \rho.x J_{\rho,y}}{\sum_\rho J_{\rho,y}}$$

Since $\mathbb{D}\mathsf{avg}$ simply distributes through the inners of $[\![J]\!]$, we deduce that $\mathbb{D}\mathsf{avg}.[\![J]\!]$ has the exact same inners as $[\![M \cdot J]\!]$ with the exact same weights. That is, the two discrete hyper-distributions are equal.  □


APPENDIX E. PROPERTIES OF THE REFINEMENT ORDER ($\sqsubseteq$)


## E.1. **Abstract *HMM*'s are ($\sqsubseteq$)-monotonic.** [§8.1]

Super-linearity (Lem. 8.3) is equivalently ($\sqsubseteq$)-monotonicity of the Kleisli-extension $h^\dagger$ of any $h{:}\mathbb{H}\mathcal{X}$; that is, it is equivalent to the more general $\Delta_1{\sqsubseteq}\Delta_2 \Rightarrow h^\dagger.\Delta_1 \sqsubseteq h^\dagger.\Delta_2$. Assuming ($\sqsubseteq$)-monotonicity and recalling that $[\cdot]$ is the point distribution, we have trivially the inequality $[\pi_1]_p{+}[\pi_2] \sqsubseteq [\pi_{1p}{+}\pi_2]$ and so

$$\begin{array}{lll}
& h.\pi_{1\ p}{+}\ h.\pi_2 & \\
= & h^\dagger.[\pi_1]_{\ p}{+}\ h^\dagger.[\pi_2] & \text{``defn. } h^\dagger\text{''} \\
= & h^\dagger.([\pi_1]_p{+}[\pi_2]) & \text{``}h^\dagger \text{ linear''} \\
& & \\
\sqsubseteq & h^\dagger.[\pi_{1p}{+}\pi_2] & \text{``}[\pi_1]_p{+}[\pi_1] \sqsubseteq [\pi_{1p}{+}\pi_2]; \\
& & \text{assumption that } h^\dagger \text{ is monotonic''} \\
& & \\
= & h.(\pi_{1p}{+}\pi_2) \ . & \text{``defn. } h^\dagger\text{''}
\end{array}$$

For the other direction (sketch), in the discrete case we note that a proof of $\Delta_1{\sqsubseteq}\Delta_2$ can be broken down into a succession of column-merges (in the matrix representation), each of them being of the form "replace $[\pi_1]_p{+}[\pi_2]$ by $[\pi_{1p}{+}\pi_2]$".

### E.2. Composition of abstract *HMM*'s respects the refinement order.          [§8.1]

We show that sequential composition of abstract *HMM*'s respects the refinement order ($\sqsubseteq$) on both sides, i.e. that for $h, h_{1,2} : \mathbb{H}\mathcal{X}$ we have both

$$h_1 \sqsubseteq h \quad \Rightarrow \quad h_1; h_2 \sqsubseteq h; h_2 \tag{E.1}$$

$$\text{and} \quad h_2 \sqsubseteq h \quad \Rightarrow \quad h_1; h_2 \sqsubseteq h_1; h \ . \tag{E.2}$$

Although this can be argued directly in terms of abstract *HMM*'s, it is easier if we use the *UM*'s defined later (§9). For (E.1) we have

$$\quad\quad h_1 \sqsubseteq h$$
iff $\quad\quad \text{wp}.h_1 \le \text{wp}.h$                                             "§E.3 just below"
implies
$$\quad\quad \text{wp}.h_1 \circ \text{wp}.h_2 \ \le \ \text{wp}.h \circ \text{wp}.h_2$$
iff $\quad\quad \text{wp}.(h_1; h_2) \le \text{wp}.(h_1; h)$                             "Cor. 9.9 in §H"
iff $\quad\quad h_1; h_2 \sqsubseteq h_1; h$ .                                                "§E.3"

And for (E.2) we have

$$\quad\quad h_2 \sqsubseteq h$$
iff $\quad\quad \text{wp}.h_2 \le \text{wp}.h$                                             "§E.3 just below"

implies                                              "wp.$h_1$ is ($\le$)-monotonic, Lem. 9.5(2)"
$$\quad\quad \text{wp}.h_1 \circ \text{wp}.h_2 \ \le \ \text{wp}.h_1 \circ \text{wp}.h$$

iff $\quad\quad \text{wp}.(h_1; h_2) \le \text{wp}.(h_1; h)$                             "Cor. 9.9 in §H"
iff $\quad\quad h_1; h_2 \sqsubseteq h_1; h$ .                                                "§E.3"

### E.3. Refinement of transformers.                        [§9.4]

Here we prove the correspondence between the forwards- and the backwards manifestations of refinement ($\sqsubseteq$), i.e. that we have

$$h_1 \sqsubseteq h_2 \quad \text{iff} \quad \text{wp}.h_1 \le \text{wp}.h_2 \ ,$$

where on the *rhs* we have extended ($\le$) pointwise, i.e. meaning $\text{wp}.h_1.u.\pi \le \text{wp}.h_1.u.\pi$ for all $u : \mathbb{U}\mathcal{X}$ and $\pi : \mathbb{D}\mathcal{X}$. We reason

$$\quad\quad h_1 \sqsubseteq h_2$$
iff                                                "pointwise extension ($\sqsubseteq$)"
$$\quad\quad h_1.\pi \sqsubseteq h_2.\pi \quad \text{for all } \pi : \mathbb{D}\mathcal{X}$$

iff                                        "Lem. 9.2, soundness and completeness"
$$\quad\quad \mathcal{E}_{h_1.\pi}\, u \le \mathcal{E}_{h_2.\pi}\, u \quad \text{for all } \pi : \mathbb{D}\mathcal{X} \text{ and all } u : \mathbb{U}\mathcal{X}$$

iff                                                  "defn. wp.()"
$$\quad\quad \text{wp}.h_1.u.\pi \le \text{wp}.h_2.u.\pi \quad \begin{array}{l} \text{for all } \pi : \mathbb{D}\mathcal{X} \\ \text{and for all } u : \mathbb{U}\mathcal{X} \end{array}$$

iff $\quad\quad \text{wp}.h_1 \le \text{wp}.h_2$ .                                           "pointwise extension"

### E.4. **Composition respects transformer refinement.** [§9.4]

For $t_{1,2}\colon\mathbb{U}\mathcal{X}$ we have defined $t_1\sqsubseteq t_2$ to be simply that $t_1.u{\leq}t_2.u$ for all $u\colon\mathbb{U}\mathcal{X}$. Here we show that functional composition of transformers respects that refinement order ($\sqsubseteq$) on both sides, i.e. that for $t,t_{1,2}\colon\mathbb{T}\mathcal{X}$ we have both $t_1{\sqsubseteq}t \Rightarrow t_1{\circ}t_2 \sqsubseteq t{\circ}t_2$ and $t_2 \sqsubseteq t \Rightarrow t_1{\circ}t_2 \sqsubseteq t_1{\circ}t$.

In fact it is trivial from the property (imposed by $\mathbb{T}\mathcal{X}$) that transformers are ($\leq$)-monotonic, that is Lem. 9.5(2).

### E.5. **Soundness and completeness: Lem. 9.2.** [§9]

We mention soundness and completeness in this paper because it provides an important justification for our definition and use of the general uncertainty measures and, in particular, their transformers.

The <u>soundness</u> part of Lem. 9.2 is related to the Data-Processing Inequality, the *DPI* [13] , which concerns two channels $C\colon\mathcal{X}{\dashrightarrow}\mathcal{Y}$ and $R\colon\mathcal{Y}{\dashrightarrow}\mathcal{Z}$. (Note that the channel $R$ here takes the *observations* $\mathcal{Y}$ of Channel $C$ as its input. Our *HMM*'s do not take observations as input.)

Informally stated, the *cascade* of $C$ and $R$ is the channel given by the matrix multiplication $C{\cdot}R$, and the *DPI* states that the information leakage from $C{\cdot}R$ cannot be more than the leakage from $C$ alone: adding another child to the game "Chinese Whispers" cannot make the eventual output less ridiculous.

We call this *soundness* because it states that a no-less-secure hyper wrt. our uncertainty refinement order indeed cannot be less uncertain when tested with *any* uncertainty measure. This result is proved in in [6,8].

The <u>completeness</u> part of Lem. 9.2 is related to the "Coriaceous Conjecture" partially proved in [8], which became the Coriaceous Property (*CP*) in [6] when its proof, for channels, was presented in complete form based on McIver's earlier, complete proof in [4] for hypers. [22] In [6] terms, the *CP* is that if there *is no R* such that $C_1{\cdot}R{=}C_2$ then there is a "gain function" (for us here, a loss function, which determines a special form of uncertainty measure in our terms) that is witness to the non-existence of such an $R$. The importance of the *CP* for quantitative information-flow security was explained in [8], and it was proved there to hold for many interesting special cases of $C_{1,2}$. But not for all of them.

The *CP* was proved to extend beyond the discrete case, to proper measure spaces, in [5].

## Appendix F. Proof of Lem. 9.4 [§9.2]

This technical lemma assures the well definedness of our dual space: we have defined our uncertainty measures $\mathbb{U}\mathcal{X}$ as functions in $\mathbb{D}\mathcal{X}{\to}\mathbb{R}^{\geq}$ with certain properties; and we have stated that wp.$h.u$ is also an uncertainty measure. Thus we must show that wp.$h.u \in \mathbb{D}\mathcal{X}{\to}\mathbb{R}^{\geq}$ and that it satisfies the properties for membership of $\mathbb{U}\mathcal{X}$.

*Lemma 9.4: Well-definedness of Def. 9.3*  If $h\colon\mathbb{H}\mathcal{X}$ is an abstract *HMM* and $u\colon\mathbb{U}\mathcal{X}$ is a *UM*, then wp.$h.u$ is in $\mathbb{U}\mathcal{X}$.

*Proof.* Since $h{\in}\mathbb{H}\mathcal{X}$, we have $h\colon\mathbb{D}\mathcal{X}{\to}\mathbb{D}^2\mathcal{X}$ satisfying Lems. 8.1,8.3. We must show for any $u\colon\mathbb{U}\mathcal{X}$ and $\delta\colon\mathbb{D}\mathcal{X}$ that $\delta{\mapsto}\mathcal{E}_{h.\delta}\,u$ is in $\mathbb{U}\mathcal{X}$, i.e. that it is in $\mathbb{D}\mathcal{X} \to \mathbb{R}^{\geq}$, is concave and is continuous (Def. 9.1).

Membership of $\delta{\mapsto}\mathcal{E}_{h.\delta}\,u$ in $\mathbb{D}\mathcal{X} \to \mathbb{R}^{\geq}$ is trivial.

---

[22]Geoffrey Smith has since told us that it follows from a result of Blackwell [39].

<u>For concavity:</u>   Because $u \in \mathbb{U}\mathcal{X}$ we know it is itself concave; and we have that $h$ satisfies the properties in Def. 8.4. We now reason

$$\text{wp}.h.u.(\pi_1 \,_p{+}\, \pi_2)$$
$$= \quad \mathcal{E}_{h.(\pi_{1p}+\pi_2)}\, u \qquad\qquad\qquad\qquad\qquad\qquad\qquad \text{``Def. 9.3''}$$

$$\geq \qquad\qquad\qquad\qquad\qquad\qquad\quad \text{`` Def. 8.4, hence } (h.\pi_1)_p{+}(h.\pi_2) \sqsubseteq h.(\pi_{1p}{+}\pi_2)$$
$$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad u \text{ concave, hence Lem. 6.4 (non-strict) applies ''}$$

$$\mathcal{E}_{(h.\pi_1)_p+(h.\pi_2)}\, u$$

$$= \quad \mathcal{E}_{h.\pi_1}\, u \,_p{+}\, \mathcal{E}_{h.\pi_2}\, u \qquad\qquad\qquad\qquad\qquad\qquad\quad \text{``}\mathcal{E}\text{ is linear''}$$
$$= \quad \text{wp}.h.u.\pi_1 \,_p{+}\, \text{wp}.h.u.\pi_2 \;, \qquad\qquad\qquad\qquad\quad \text{``Def. 9.3''}$$

as required.

    <u>For continuity:</u>   We must show that $\text{wp}.h.u$ is continuous, given that both $u, h$ are themselves continuous. Because $h$ itself is continuous, we need only show that in general the function $\Delta \mapsto \mathcal{E}_\Delta\, u$ is continuous wrt. Kantorovich on the left, in $\Delta$ for fixed continuous $u$. This follows from the fact that $\mathbb{D}\mathcal{X}$ is a compact metric space so that the Kantorovich metric metrizes the weak topology on $\mathbb{D}^2\mathcal{X}$. That is, $\mathcal{E}_{\Delta_n}\, u$ converges to $\mathcal{E}_\Delta\, u$ for every continuous function $u{:}\,\mathbb{D}\mathcal{X} \to \mathbb{R}^{\geq}$ iff $\Delta_n$ converges to $\Delta$ wrt. Kantorovich metric. $\qquad\square$

## Appendix G. Extension of transformers      [§9.4]

The core ingredient in the proof of this theorem is the Riesz Representation Theorem for linear functionals (linear maps from a normed vector space to $\mathbb{R}$). A difficulty however originates from the fact that the representation theorem is stated on the space of all continuous functions $\mathbb{C}\mathcal{X}$ (defined below), but our linear function $t$ is defined only from the subspace $\mathbb{U}\mathcal{X}$ to itself.

**Definition G.1** (Space of continuous functions)**.** We define $\mathbb{C}\mathcal{X}$ to be the set of all *continuous functions* from $\mathbb{D}\mathcal{X}$ (with the Kantorovich metric) to $\mathbb{R}$ (with the ordinary metric). This set is endowed with the *uniform metric* $\| \cdot - \cdot \|_\infty$, defined

$$\|u_1 - u_2\|_\infty \quad := \quad \sup_{\delta:\, \mathbb{D}\mathcal{X}} |u_1.\delta - u_2.\delta| \;, \qquad\qquad\qquad (\text{G.1})$$
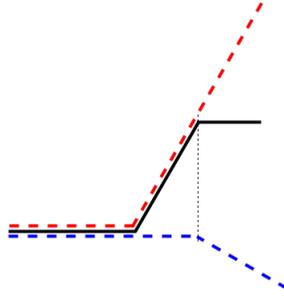
that turns $\mathbb{C}\mathcal{X}$ into a complete metric space.

    Yet $\mathbb{U}\mathcal{X}$ is a sub-metric space of $\mathbb{C}\mathcal{X}$ under the uniform metric $\| \cdot - \cdot \|_\infty$. More importantly, we prove that the vector space generated by $\mathbb{U}\mathcal{X}$ is dense in $\mathbb{C}\mathcal{X}$. (See Lem. G.2 and Fig. 7.) This is essential to ensure that if $t$ extends to a continuous linear function over $\mathbb{C}\mathcal{X}$, then such an extension is necessarily unique. We will show in Thm. G.4 that such an extension always exists.

**Lemma G.2** (Concave density)**.** *The vector space generated by $\mathbb{U}\mathcal{X}$ is dense in $\mathbb{C}\mathcal{X}$ wrt. $\| \cdot - \cdot \|_\infty$.*

*Proof.* This result essentially follows from [40, Pro. 2.2]. We give the proof here for completeness.

    Let $\langle \mathbb{U}\mathcal{X} \rangle$ be the set of functions that can be written as the difference of two positive concave functions from $\mathbb{D}\mathcal{X}$ to $\mathbb{R}$. Then $\langle \mathbb{U}\mathcal{X} \rangle$ coincides with the real vector space generated

Every continuous piecewise-linear function can be expressed as the sum of finitely many convex and concave functions using the construction shown in this figure. This provides a geometrical view of the Concave Density Lem. G.2. By summing up concave and convex functions, we get non concave (resp. convex) functions. In fact, the density theorem shows that any arbitrarily shaped function can be approximated using finite sums and products of scaled concave functions.

FIGURE 7. The sum of convex (red, upper) and concave (blue, lower) functions gives a zig-zag (black, middle).

by $\mathbb{U}\mathcal{X}$ (by grouping positively and negatively weighted components). Equivalently, every function in $\langle\mathbb{U}\mathcal{X}\rangle$ is the difference of two positive continuous convex functions: if $f = u_1 - u_2$ for $u_{1,2}{:}\,\mathbb{U}\mathcal{X}$, then
$$f = (-u_2 - c) - (-u_1 - c)$$
where $c = \min(\inf_{\delta:\,\mathbb{D}\mathcal{X}} -u_1.\delta, \inf_{\delta:\,\mathbb{D}\mathcal{X}} -u_2.\delta)$. The constant $c$ is finite because $\mathbb{D}\mathcal{X}$ is compact. The functions $-u_{1,2} - c$ are positive, continuous and convex functions.

Now let us apply the Stone-Weierstrass Density Theorem [25, Thm. 5] on $\langle\mathbb{U}\mathcal{X}\rangle$ which is a subset of $\mathbb{C}\mathcal{X}$.

To do that we need first to show that $\langle\mathbb{U}\mathcal{X}\rangle$ is an algebra (i.e. has a zero and unit, is closed under scalar multiplication and addition and pointwise multiplication of $f$'s). In addition $\langle\mathbb{U}\mathcal{X}\rangle$ must "vanish nowhere" on $\mathbb{D}\mathcal{X}$ and "separate points". (See below for explanations of those properties.)

$\underline{\langle\mathbb{U}\mathcal{X}\rangle \text{ is an algebra}}$: Since $\langle\mathbb{U}\mathcal{X}\rangle$ is a vector space, the constant functions $\mathbf{0}, \mathbf{1}$ (identically 0 and 1 resp.) and the functions $cf, f + g$ are in $\langle\mathbb{U}\mathcal{X}\rangle$ for every $c{:}\,\mathbb{R}$ and $f, g{:}\,\langle\mathbb{U}\mathcal{X}\rangle$.

Let $f, g{:}\,\langle\mathbb{U}\mathcal{X}\rangle$ be such that $f = u_1 - u_2$ and $g = v_1 - v_2$ where $u_{1,2}, v_{1,2}$ are positive continuous convex functions. Notice that
$$f^2 \quad = \quad 2(u_1^2 + u_2^2) - (u_1 + u_2)^2$$
where $u_{1,2}^2$ and $(u_1 + u_2)^2$ are positive convex functions (because the square of a non-negative convex function is convex). That is, we have $f^2 \in \langle\mathbb{U}\mathcal{X}\rangle$. Now
$$fg \quad = \quad (f + g)^2 - (f^2 + g^2)$$
and thus $fg \in \langle\mathbb{U}\mathcal{X}\rangle$, because we have just shown that all of $(f + g)^2, f^2, g^2$ are in $\langle\mathbb{U}\mathcal{X}\rangle$.

$\underline{\langle\mathbb{U}\mathcal{X}\rangle \text{ vanishes nowhere}}$: We must show that for each $\delta{:}\,\mathbb{D}\mathcal{X}$ there is some $f{:}\,\langle\mathbb{U}\mathcal{X}\rangle$ such that $f\delta \neq 0$. But this is immediate since $\mathbf{1}.\delta \neq 0$ for every $\delta{:}\,\mathbb{D}\mathcal{X}$ and $\mathbf{1} \in \langle\mathbb{U}\mathcal{X}\rangle$.

$\langle \mathbb{U}\mathcal{X} \rangle$ separates points: We must show that for every pair $\delta \neq \delta'$ there is some $f \in \langle \mathbb{U}\mathcal{X} \rangle$ such that $f.\delta \neq f.\delta'$. We argue as follows.

Given $\delta \colon \mathbb{D}\mathcal{X}$ (fixed) and define $f^\delta.\delta' := d_K(\delta', \delta)$ for $\delta' \colon \mathbb{D}\mathcal{X}$. Observe that for $\delta' \neq \delta$ we have that $0 = f^\delta.\delta < d_K(\delta', \delta) = f^\delta.\delta'$. Thus it suffices to show that $f^\delta \in \langle \mathbb{U}\mathcal{X} \rangle$, and we continue as follows.

For every $\delta_{1,2} \colon \mathbb{D}\mathcal{X}$, we have

$$
\begin{array}{lll}
& f^\delta.(\delta_{1p} + \delta_2) & \\
= & d_K((\delta_{1p} + \delta_2), \delta) & \text{``}d_K \text{ is Kantorovich distance; Definition of } f^\delta\text{''} \\
\leq & d_K(\delta_1, \delta)_p + d_K(\delta_2, \delta) & \text{``}d_K(\delta', \delta) \text{ is convex, for fixed } \delta\text{''} \\
= & f^\delta.\delta_{1p} + f^\delta.\delta_2 & \text{``Definition of } f^\delta\text{''}
\end{array}
$$

That is $\mathbf{1} - f^\delta \in \mathbb{U}\mathcal{X}$ and thus $\langle \mathbb{U}\mathcal{X} \rangle$ separates points.

By the Stone-Weierstrass Theorem [25, Thm. 5], we have $\langle \mathbb{U}\mathcal{X} \rangle$ is dense in $\mathbb{C}\mathcal{X}$. □

The extension of a transformer $t \colon \mathbb{U}\mathcal{X} \to \mathbb{U}\mathcal{X}$ to a continuous linear function from $\mathbb{C}\mathcal{X}$ to itself is done in two stages. Firstly, $t$ is extended *linearly* to a continuous linear function $t' \colon \langle \mathbb{U}\mathcal{X} \rangle \to \mathbb{C}\mathcal{X}$. This step is justified in Thm. G.4. Secondly, $t'$ is extended *continuously* to a continuous linear function $\tilde{t} \colon \mathbb{C}\mathcal{X} \to \mathbb{C}\mathcal{X}$. This step uses the density proven in Lem. G.2 and is shown in Lem. G.3 below.

**Lemma G.3** (Extension from $\langle \mathbb{U}\mathcal{X} \rangle$ to $\mathbb{C}\mathcal{X}$). *Every continuous linear function $t$ from $\langle \mathbb{U}\mathcal{X} \rangle$ to $\mathbb{C}\mathcal{X}$ extends uniquely to a continuous linear function $\tilde{t}$ from $\mathbb{C}\mathcal{X}$ to itself.*

*Proof.* This result follows from the Continuous Linear Extension Theorem [41, Ch. 4 Thm. 10.1].

All we need to show is that the (Cauchy) completion of $\langle \mathbb{U}\mathcal{X} \rangle$ is $\mathbb{C}\mathcal{X}$, which follows from the fact that $\langle \mathbb{U}\mathcal{X} \rangle$ is dense in $\mathbb{C}\mathcal{X}$ (Lem. G.2) and that $\mathbb{C}\mathcal{X}$ is a complete normed vector space when endowed with the uniform norm $\|f\|_\infty := \|f - \mathbf{0}\|_\infty$. □

**Theorem G.4** (Extension from $\mathbb{U}\mathcal{X}$ to $\mathbb{C}\mathcal{X}$). *Every transformer extends uniquely to a positive continuous linear function from $\mathbb{C}\mathcal{X}$ to itself.*

*Proof.* Let $t \colon \mathbb{T}\mathcal{X}$ be a transformer. It suffices to prove that $t$ has a positive continuous extension $t'$ on the sub-vector space $\langle \mathbb{U}\mathcal{X} \rangle$. If such a $t'$ exists then a unique extension $\tilde{t} \colon \mathbb{C}\mathcal{X} \to \mathbb{C}\mathcal{X}$, which is positive [23] and continuous, can be deduced using Lem. G.3.

Let $f \colon \langle \mathbb{U}\mathcal{X} \rangle$, there exists $u_{1,2} \in \mathbb{U}\mathcal{X}$ such that $f = u_1 - u_2$. We define $t'.f = t.u_1 - t.u_2$.

$\underline{t' \text{ is well-defined}}$:   We must show that $t'.f$ is independent of how $f$ is written as the difference of two uncertainty measures. Firstly, notice that if $u_1 - u_2 \in \mathbb{U}\mathcal{X}$ for some $u_{1,2} \in \mathbb{U}\mathcal{X}$ then $t.(u_1 - u_2) = t.u_1 - t.u_2$. Secondly, let $f = u_1 - u_2 = v_1 - v_2$. Then $(u_1 + v_2) - (u_2 + v_1) = \mathbf{0}$, which is in $\mathbb{U}\mathcal{X}$. Therefore, we have $t.(u_1 + v_2) - t.(u_2 + v_1) = \mathbf{0}$, and that implies $t.u_1 - t.u_2 = t.v_1 - t.v_2$ by linearity of $t$.

$\underline{t' \text{ is linear and unique}}$:   Linearity is clear and it implies uniqueness of the extension $t'$ over $\langle \mathbb{U}\mathcal{X} \rangle$.

$\underline{t' \text{ is 1-Lipschitz}}$:   Let $f, g \colon \langle \mathbb{U}\mathcal{X} \rangle$ be such that we have $f = u_1 - u_2$ and $g = v_1 - v_2$. Then

---

[23]For the positiveness of the continuous extension, if $f$ is a positive continuous function that is the uniform limit of a sequence of $f_n$'s in $\langle \mathbb{U}\mathcal{X} \rangle$, then the sequence of positive continuous functions $\max(\mathbf{0}, f_n) \in \langle \mathbb{U}\mathcal{X} \rangle$ also converges to $f$ wrt. the uniform metric. The reason is $|f.\delta - \max(0, f_n.\delta)| \leq |f.\delta - f_n.\delta|$, for every $\delta \colon \mathbb{D}\mathcal{X}$ and positive $f$. Thus $t.f$ has to be positive.

$$\|t'.f - t'.g\|_\infty$$
$$= \quad \|(t.u_1 - t.u_2) - (t.v_1 - t.v_2)\|_\infty \qquad\qquad \text{“Definition of } t'\text{”}$$
$$= \quad \|t.(u_1 + v_2) - t.(v_1 + u_2)\|_\infty \qquad\qquad \text{“}t \text{ is linear, } u_i + v_j \in \mathbb{U}\mathcal{X}\text{”}$$
$$\leq \quad \|(u_1 + v_2) - (v_1 + u_2)\|_\infty \qquad\qquad \text{“}t \text{ is 1-Lipschitz”}$$
$$= \quad \|f - g\|_\infty \; . \qquad\qquad \text{“Definition of } f, g\text{”}$$

Therefore, $t'$ is also continuous.

$t'$ is positive: (i.e. it maps non-negative functions to non-negative functions). This follows from monotonicity of $t$.

By Lem. G.3, the extension $t'$ further extends into a continuous positive linear function $\tilde{t} \colon \mathbb{C}\mathcal{X} \to \mathbb{C}\mathcal{X}$ with $\tilde{t}.u = t.u$ for every $u \colon \mathbb{U}\mathcal{X}$. $\qquad\square$

## APPENDIX H. PROOF OF COR. 9.9 [§9.4]

This proof is made easier by operating in a slightly more general space than $\mathbb{H}\mathcal{X}$, i.e. the measurable subset of $\mathbb{D}\mathcal{X} \to \mathbb{D}^2\mathcal{X}$, not taking advantage of the stronger conditions that characterise $\mathbb{H}\mathcal{X}$ within it. In this section only we write $\overline{\mathrm{wp}}.$ for the function defined as at Def. 9.3 but over the larger space.

**Lemma H.1** (Transformer composition). *For any (measurable) $h_{1,2} \colon \mathbb{D}\mathcal{X} \to \mathbb{D}^2\mathcal{X}$ we have that $\overline{wp}.(h_1; h_2) = \overline{wp}.h_1 \circ \overline{wp}.h_2$.*

*Proof.* $\quad \overline{\mathrm{wp}}.(h_1; h_2).u.\pi$
$$= \quad \mathcal{E}_{(h_1;h_2).\pi}\, u \qquad\qquad \text{“Def. 9.3 extended to } \overline{\mathrm{wp}}.()\text{”}$$
$$= \quad \mathcal{E}_{\mathsf{avg}.(\mathbb{D}h_2.(h_1.\pi))}\, u \qquad\qquad \text{“}h_1; h_2 \text{ is Kleisli composition”}$$

$$= \quad \mathcal{E}_{\mathbb{D}h_2.(h_1.\pi)}\,(\lambda\Delta \cdot \mathcal{E}_\Delta\, u) \qquad \begin{array}{l}\text{“}\mathcal{E}_{\mathsf{avg}.\underline{\Delta}}\, u = \mathcal{E}_{\underline{\Delta}}\,(\lambda\Delta \cdot \mathcal{E}_\Delta\, u) \text{ from (†) below} \\ \lambda \text{ is lambda-abstraction”}\end{array}$$

$$= \quad \mathcal{E}_{h_1.\pi}\,((\lambda\Delta \cdot \mathcal{E}_\Delta\, u) \circ h_2) \qquad\qquad \text{“}\mathcal{E}_{\mathbb{D}h_2.\Delta}\, F = \mathcal{E}_\Delta\,(F \circ h_2) \text{ from (‡) below”}$$
$$= \quad \mathcal{E}_{h_1.\pi}\,(\lambda\pi' \cdot ((\lambda\Delta \cdot \mathcal{E}_\Delta\, u) \circ h_2).\pi') \qquad\qquad \text{“make } \pi' \text{ explicit”}$$
$$= \quad \mathcal{E}_{h_1.\pi}\,(\lambda\pi' \cdot \mathcal{E}_{h_2.\pi'}\, u) \qquad\qquad \text{“}\Delta := h_2.\pi'\text{”}$$
$$= \quad \overline{\mathrm{wp}}.h_1.(\lambda\pi' \cdot \mathcal{E}_{h_2.\pi'}\, u).\pi \qquad\qquad \text{“Def. 9.3”}$$
$$= \quad \overline{\mathrm{wp}}.h_1.(\overline{\mathrm{wp}}.h_2.u).\pi \; . \qquad\qquad \text{“Def. 9.3”}$$

The identities (†) and (‡) were proven by Giry ([2, Sec. 3 p.70]). In (‡), $F$ maps every hyper $\Delta$ to $\mathcal{E}_\Delta\, u$. $\qquad\square$

Remarkably, it is quite easy to show that wp.() is an injection over all of $\mathbb{D}\mathcal{X} \to \mathbb{D}^2\mathcal{X}$.

**Lemma H.2** (wp.() is an injection on $\mathbb{D}\mathcal{X} \to \mathbb{D}^2\mathcal{X}$). *If $wp.h_1 = wp.h_2$ for some (measurable) $h_{1,2} \colon \mathbb{D}\mathcal{X} \to \mathbb{D}^2\mathcal{X}$, then $h_1 = h_2$.*

*Proof.* We reason

$$\quad h_1 \neq h_2$$
$$\Rightarrow \quad h_1.\pi \neq h_2.\pi \qquad\qquad \text{“for some } \pi \colon \mathbb{D}\mathcal{X}\text{”}$$
$$\Rightarrow \quad h_1.\pi \not\sqsubseteq h_2.\pi \qquad\qquad \text{“}wlog; \; (\sqsubseteq)\text{-antisymmetry from §6”}$$
$$\Rightarrow \quad \mathcal{E}_{h_2.\pi}\, u < \mathcal{E}_{h_1.\pi}\, u \qquad\qquad \text{“Lem. 9.2 completeness (}Coriaceous\text{),for some } u \colon \mathbb{U}\mathcal{X}\text{”}$$
$$\Leftrightarrow \quad \mathrm{wp}.h_2.u.\pi < \mathrm{wp}.h_1.u.\pi \qquad\qquad \text{“Def. 9.3”}$$
$$\Rightarrow \quad \mathrm{wp}.h_1 \neq \mathrm{wp}.h_2 \; .$$

$\square$

Our next step is to use Thm. 9.8 to show that indeed $h_1; h_2 \in \mathbb{H}\mathcal{X}$, so that $\overline{\text{wp}}.$ can be replaced by wp.() in Lem. H.1 just above. We have

**Lemma H.3** ($\mathbb{H}\mathcal{X}$ closed under composition)**.** *For $h_{1,2}: \mathbb{H}\mathcal{X}$ we have $h_1; h_2 \in \mathbb{H}\mathcal{X}$.*

*Proof.* If $h_{1,2}: \mathbb{H}\mathcal{X}$ then wp.$h_{1,2} \in \mathbb{T}\mathcal{X}$ from Lems. 9.5,9.6; and since those properties are closed under composition, we have that wp.$h_1 \circ$ wp.$h_2 \in \mathbb{T}\mathcal{X}$ as well.

From Thm. 9.8 there is then a unique $h: \mathbb{H}\mathcal{X}$ such that wp.$h$ = wp.$h_1 \circ$ wp.$h_2$; but examination of Lem. H.2 shows membership of $\mathbb{H}\mathcal{X}$ is not necessary for that uniqueness: it applies to the whole of (measurable) $\mathbb{D}\mathcal{X} \rightarrow \mathbb{D}^2\mathcal{X}$. That is, there no other measurable $h$ in all of $\mathbb{D}\mathcal{X} \rightarrow \mathbb{D}^2\mathcal{X}$ such that $\overline{\text{wp}}.h = t$.

From Lem. H.1 we know that $\overline{\text{wp}}.(h_1; h_2) = t$, and so we must have $h_1; h_2 = h \in \mathbb{H}\mathcal{X}$. $\square$

Thus we can conclude

*Corollary 9.9: Transformer composition* For any $h_{1,2}: \mathbb{H}\mathcal{X}$ we have that also $h_1; h_2 \in \mathbb{H}\mathcal{X}$, and furthermore wp.$(h_1; h_2)$ = wp.$h_1 \circ$ wp.$h_2$.

*Proof.* Lemmas H.1,H.3 just above. $\square$

## Appendix I. Calculation of wp.$[\![P]\!]$ [§11.2]

In §11.2 a sample analysis was done on a very small program to show how, if the post-uncertainty is fixed, a pre-uncertainty can be calculated once and for all; and that then that pre-uncertainty can be used to investigate the security implications of a number of different priors, without having to re-analyse the program for each one.

Here we give the calculations for wp.$(\cdot)$ in §11.2. We note below however that ideally the pre-uncertainty would be calculated by *source-level* reasoning; but that is not what we do here. (See also our "more concrete aim" in §14 concerning source-level reasoning.)

Let $P$ be the program set out in Fig. 4 (and also Fig. 5 from §4.4). As usual for weakest preconditions, we work from post- to pre-. Let $u$ be the *UM* from §11.1, reflecting the circumstances of an attacker whose principal concern is whether the two bits of `xs` are the same.

Beginning with the second statement, since with transformers we work from the back towards the front, we expect informally that wp.$[\![\texttt{xs:= xs}_{1/2}\oplus\texttt{-xs}]\!].u$ is just $u$ again — since the assignment does not affect `xs[0]=xs[1]`, whichever branch is taken. Calculation confirms that: for arbitrary $\pi$ we have

wp.$[\![\texttt{xs:= xs}_{1/2}\oplus\texttt{-xs}]\!].u.\pi$

$$= \quad \mathcal{E}_{[(\,(\pi_{00}+\pi_{11})/2, (\pi_{01}+\pi_{10})/2,\ \ u}_{(\pi_{10}+\pi_{01})/2, (\pi_{11}+\pi_{00})/2\,)]} \qquad \text{``semantics of } \texttt{xs:= xs}_{1/2}\oplus\texttt{-xs''}$$

$$= \quad u.\,(\,(\pi_{00}+\pi_{11})/2, (\pi_{01}+\pi_{10})/2, \qquad\qquad \text{``expectation over point hyper''}$$
$$(\pi_{10}+\pi_{01})/2, (\pi_{11}+\pi_{00})/2\,)$$

$$= \quad \min \begin{array}{l} (\pi_{00}+\pi_{11})/2 + (\pi_{11}+\pi_{00})/2 \\ (\pi_{01}+\pi_{10})/2 + (\pi_{10}+\pi_{01})/2 \end{array} \qquad \text{``definition } u \text{ from §11.2''}$$

$$= \qquad (\pi_{00}+\pi_{11}) \ \min \ (\pi_{01}+\pi_{10})$$
$$= \qquad u \ , \qquad\qquad\qquad\qquad\qquad\qquad\text{``definition } u \text{ again''}$$

as we expected.

Continuing towards the front of the program we now calculate again for arbitrary $\pi$, but from just above able to use the same $u$ that we started with, that

$$\text{wp.}[\![\texttt{leak xs}[0]\,_{1/2}\oplus\texttt{xs}[1]]\!].u.\pi$$

$$= \qquad \mathcal{E} \underset{\substack{ \\ s_0\oplus\, (0,\pi_{01}/2s_1,\pi_{10}/2s_1,\pi_{11}/s_1)}}{(\pi_{00}/s_0,\pi_{01}/2s_0,\pi_{10}/2s_0,0)}\ u \qquad \begin{array}{r}\text{`` semantics of } \texttt{print xs}[0]\,_{1/2}\oplus\texttt{xs}[1]\\ \text{define } s_0\!:=\!\pi_{00}\!+\!(\pi_{01}\!+\!\pi_{10})/2\\ s_1\!:=\!(\pi_{01}\!+\!\pi_{10})/2+\pi_{11}\text{ ''}\end{array}$$

$$= \qquad \begin{array}{l} u.(\pi_{00}/s_0,\pi_{01}/2s_0,\pi_{10}/2s_0,0)\quad\text{``}\mathcal{E}\text{ linear, applied to two-point hyper (Def. 2.5)''}\\ {}_{s_0\oplus}\ \ u.(0,\pi_{01}/2s_1,\pi_{10}/2s_1,\pi_{11}/s_1)\end{array}$$

$$= \qquad \pi_{00}\min(\pi_{01}+\pi_{10})/2 + (\pi_{01}+\pi_{10})/2\min\pi_{11} \ , \qquad\qquad \begin{array}{r}\text{``definition } u\\ \text{from previous calculation''}\end{array}$$

as claimed in §11.2.

We stress that calculating wp.() this way for any but the smallest programs is *not practical at all*. For a practical calculus, instead the formulation of uncertainties as loss functions would be used to write them as expressions at the source level, i.e. over program variables, and then using formal manipulations in a quantitative program logic (extending e.g. [21, 22]).

The issue of source-level reasoning is discussed further in the conclusion §14.


## Appendix J. Using loss functions to characterise pure channels

With uncertainty transformers, we can be more precise about the properties satisfied by pure-(abstract) channel *HMM*'s specifically. As with markovs the mechanism by which information is released is independent of the (probability) values associated with the prior; in fact it only depends on the underlying state value, that is $\mathcal{X}$. This property can be described neatly in terms of a "multiplicative property" on transformers which, in addition, provides a characterisation of transformers which correspond to channels. We begin with a motivating example.

Take $\mathcal{X}=\{0,1\}$. It's easy to construct an $h\colon\mathbb{H}\mathcal{X}$ with the property that for all $\pi\colon\mathbb{D}\mathcal{X}$ we have $\textsf{avg.}(h.\pi)=\pi$, which is to say that its markov is the identity, but it is still not a pure channel: we simply "cheat" by using a different channel for each prior. Take for example the $\pi$-indexed channels given by the matrix

$$C^\pi \quad:=\quad \begin{pmatrix} \pi_0 & \pi_1 \\ 0 & 1 \end{pmatrix} \ .$$

The function defined $f.\pi:=[\![\pi\triangleright C^\pi]\!]$ does not satisfy $f=[\![C\!:]\!]$ for any *single* fixed $C$, and this example provides the insight for characterising pure channels: they have a simple multiplicative property, which we express using loss functions as follows.

**Definition J.1** (Multiplicativity of transformers)**.** For loss-function $l: I \to \mathcal{X} \to \mathbb{R}^{\geq}$ and $\pi: \mathbb{D}\mathcal{X}$ define a $\pi$-skewed loss function $(l \lhd \pi).i.x := l.i.x \times \pi.x$. We then say that transformer $t: \mathbb{T}\mathcal{X}$ is *multiplicative* if for any $\pi_{1,2}: \mathbb{D}\mathcal{X}$ and loss function $l$ we have $t.(U_{l \lhd \pi_1}).\pi_2 = t.(U_{l \lhd \pi_2}).\pi_1$. [24]

**Lemma J.2** (Channels are multiplicative)**.** *Let $C: \mathcal{X} \nrightarrow \mathcal{Y}$ be a channel matrix. Then $\mathrm{wp}.[\![C{:}]\!]$ is multiplicative.*

*Proof.* This follows because the identity transformer is multiplicative, i.e. $(U_{l \lhd \pi_1}).\pi_2 = (U_{l \lhd \pi_2}).\pi_1$, and that wp.() applied to a pure channel maps any given loss function to a sum of loss functions "scaled" by the columns. $\qquad\square$

The following fact shows that this multiplicative property in fact characterises channels.

**Lemma J.3.** *Let $f: \mathbb{D}\mathcal{X} \to \mathbb{D}^2\mathcal{X}$ be such that $f.\pi$ has finite support for every $\pi: \mathbb{D}\mathcal{X}$; assume it satisfies the pure-channel property from §B.2; and assume that wp.$f$ is multiplicative as just above. Then there is some set of observations $\mathcal{Y}$ and channel $C: \mathcal{X} \nrightarrow \mathcal{Y}$ such that $f = [\![C{:}]\!]$.*

*Proof.* Let $N$ be the size of $\mathcal{X}$ and let $\upsilon$ be the uniform distribution on $\mathcal{X}$. [25] Define $\Delta := f.\upsilon$ and let $\mathcal{Y}$ be the support of $\Delta$, a finite set of distributions that will be used as column indices. Then define $C: \mathcal{X} \nrightarrow \mathcal{Y}$ by

$$C_{x,y} \quad := \quad N \times \Delta.y.x \;,$$

so that $f.\upsilon = \Delta = [\![C{:}]\!].\upsilon$. We now show that in fact $f.\pi = [\![C{:}]\!].\pi$ for all $\pi: \mathbb{D}\mathcal{X}$.

We have for any loss function $l$ that

$$
\begin{aligned}
&\mathcal{E}_{f.\pi}\, U_l \\
=\;& \mathrm{wp}.f.U_l.\pi \\
=\;& \mathrm{wp}.f.U_{l' \lhd \upsilon}.\pi && \text{``define } l' := N \times l\text{''} \\
=\;& \mathrm{wp}.f.U_{l' \lhd \pi}.\upsilon && \text{``assumption wp.}f \text{ multiplicative''} \\
=\;& \mathcal{E}_{f.\upsilon}\,(U_{l' \lhd \pi}) \\
=\;& \mathcal{E}_{[\![C:]\!].\upsilon}\,(U_{l' \lhd \pi}) && \text{``defn. } C\text{''} \\
=\;& \mathcal{E}_{[\![C:]\!].\pi}\, U_l \;, && \text{``reverse steps above;} \\
&&& \text{wp.}[\![C{:}]\!] \text{ multiplicative''}
\end{aligned}
$$

so $f.\pi = [\![C_\Delta{:}]\!].\pi$ since hypers are determined by loss functions [4,6], thus $f = [\![C_\Delta{:}]\!]$ because $\pi$ was arbitrary. $\qquad\square$

---

[24]This notation is by analogy with $\pi \rhd C$ that "multiplies $\pi$ in" from the $\mathcal{X}$ side of a matrix; in $C \lhd \pi$ the $\pi$ is multiplied in from the other side.

[25]It is *upsilon* for "uniform".