



**HAL**  
open science

# Refinement Metrics for Quantitative Information Flow

Konstantinos Chatzikokolakis, Geoffrey Smith

► **To cite this version:**

Konstantinos Chatzikokolakis, Geoffrey Smith. Refinement Metrics for Quantitative Information Flow. Mário S. Alvim; Kostas Chatzikokolakis; Carlos Olarte; Frank Valencia. The Art of Modelling Computational Systems: A Journey from Logic and Concurrency to Security and Privacy. Essays Dedicated to Catuscia Palamidessi on the Occasion of Her 60th Birthday., 11760, Springer, pp.397-416, 2019, Lecture Notes in Computer Science, 978-3-030-31174-2. 10.1007/978-3-030-31175-9\_23 . hal-02350777

**HAL Id: hal-02350777**

**<https://inria.hal.science/hal-02350777>**

Submitted on 6 Nov 2019

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Refinement Metrics for Quantitative Information Flow

Konstantinos Chatzikokolakis<sup>1</sup> and Geoffrey Smith<sup>2</sup>

<sup>1</sup> University of Athens [kostasc@di.uoa.gr](mailto:kostasc@di.uoa.gr)

<sup>2</sup> Florida International University [smithg@cis.fiu.edu](mailto:smithg@cis.fiu.edu)

**Abstract.** In Quantitative Information Flow, *refinement* ( $\sqsubseteq$ ) expresses the strong property that one channel never leaks more than another. Since two channels are then typically incomparable, here we explore a family of *refinement quasimetrics* offering greater flexibility. We show these quasimetrics let us unify refinement and capacity, we show that some of them can be computed efficiently via linear programming, and we establish upper bounds via the Earth Mover’s distance. We illustrate our techniques on the Crowds protocol.

## 1 Introduction

Completely eliminating the leakage of sensitive information by computer systems is often infeasible, making it attractive to approach the problem quantitatively. *Quantitative information flow* [3] offers a rich family of *g-leakage* measures of the amount of leakage caused by a channel  $C$  taking a secret input  $X$  to an observable output  $Y$ ; these measures are parameterized by  $X$ ’s *prior distribution*  $\pi$  and a *gain function*  $g$ , which models the adversary’s capabilities and goals. (See §2 for a brief review.)

While *g-leakage* lets us precisely measure information leakage in a rich variety of operational scenarios, we may be unsure about the appropriate prior  $\pi$  or gain function  $g$  to use, and hence about the robustness of our conclusions. Two approaches to robustness have proved fruitful: *capacity*, which is the maximum leakage over some sets of gain functions and priors, and *refinement*: channel  $A$  is *refined* by channel  $B$ , written  $A \sqsubseteq B$ , iff  $B$  never leaks more than  $A$ , regardless of the prior or gain function. Remarkably, refinement also has a structural characterization:  $A \sqsubseteq B$  iff there exists a “post-processing” channel matrix  $R$  such that  $B = AR$ . Moreover, ( $\sqsubseteq$ ) is a *partial order* on abstract channels.

Unfortunately, refinement is a very *partial* partial order, in that most channels are incomparable. And this “Boolean” nature of refinement is inconsistent with the spirit of quantitative information flow, which is above all motivated by the need to tolerate imperfect security. If  $A \not\sqsubseteq B$ , then we know that  $B$  can be worse than  $A$ ; but we do not know whether  $B$  is really terrible, or whether it is just slightly worse than  $A$ . This is the main issue that we address in this paper.

In mathematics, it is common to use metrics to provide a “finer”, quantified generalization of a relation. A relation has a “Boolean” nature: elements are either related or not. In the metric version, related elements have distance 0; but for non-related elements the metric tells you *how much* the relation is violated. For instance, the Euclidean distance can be seen as a quantified generalization of the equality relation ( $=$ ) on  $\mathbb{R}$  — the distance  $|x - y|$  tells us how much  $x = y$  is violated; 0 means that it is not violated

at all.<sup>3</sup> This approach is not of course limited to symmetric relations. For instance, the quasimetric

$$\mathbf{q}_{\leq}^+(x, y) := \max\{y - x, 0\},$$

can be seen as a quantified version of ( $\geq$ ); it measures how much ( $\geq$ ) is violated, and is 0 iff  $x \geq y$ .

In this paper we define a family of *refinement quasimetrics*  $\text{ref}_{\mathcal{D}, \mathcal{G}}^q$  to measure how much refinement is violated. Refinement is then the *kernel* (i.e. elements at distance 0) of these metrics: we have  $\text{ref}_{\mathcal{D}, \mathcal{G}}^q(\mathbf{A}, \mathbf{B}) = 0$  iff  $\mathbf{A} \sqsubseteq \mathbf{B}$ . Through this approach, we make the following contributions:

- We treat both additive and multiplicative leakage together via additive and multiplicative quasimetrics  $\mathbf{q}_{\leq}^+$  and  $\mathbf{q}_{\leq}^{\times}$ . (The latter crucially takes the log of the ratio, instead of just the ratio.)
- We observe that the capacity of  $\mathbf{C}$  can be expressed as  $\text{ref}_{\mathcal{D}, \mathcal{G}}^q(\mathbf{1}, \mathbf{C})$ , where  $\mathbf{1}$  is the perfect channel leaking nothing, giving us a unified way of looking at both refinement and capacity.
- In §4 we show that, for a fixed prior  $\pi$ , the additive refinement metric over all gain functions in  $\mathbb{G}^{\uparrow} \mathcal{X}$  can be computed in polynomial time via linear programming.
- In §5 we prove (again for a fixed  $\pi$ ) that the refinement metric is bounded by the Earth Mover’s distance between  $[\pi \triangleright \mathbf{A}]$  and  $[\pi \triangleright \mathbf{B}]$ , allowing it to be approximated very efficiently. We also prove a bound on the metric over all priors.
- In §6 we give a case study of our techniques on the Crowds protocol.

## 2 Preliminaries

### 2.1 Quantitative Information Flow

In *quantitative information flow* [6, 7, 9] an adversary’s prior knowledge about a secret input  $X$  drawn from a finite set  $\mathcal{X}$  is modeled as a probability distribution  $\pi$ , and a probabilistic system taking input  $X$  to observable output  $Y$  is modeled as an information-theoretic channel matrix  $\mathbf{C}$  giving the conditional probabilities  $p(y|x)$ . Assuming that the adversary knows  $\mathbf{C}$ , then each output  $y$  allows the adversary to update her knowledge about  $X$  to a posterior distribution  $p_{X|y}$ , which moreover has probability  $p(y)$  of occurring. As a result, the effect of  $\mathbf{C}$  is to map each prior distribution  $\pi$  to a *distribution on posterior distributions*, called a *hyper-distribution* and denoted  $[\pi \triangleright \mathbf{C}]$ . For example, channel

$\mathbf{C}$	$y_1$	$y_2$	$y_3$	$y_4$
$x_1$	$1/2$	$1/2$	$0$	$0$
$x_2$	$0$	$1/4$	$1/2$	$1/4$
$x_3$	$1/2$	$1/3$	$1/6$	$0$

maps  $\pi = (1/4, 1/2, 1/4)$  to hyper

$[\pi \triangleright \mathbf{C}]$	$1/4$	$1/3$	$7/24$	$1/8$
$x_1$	$1/2$	$3/8$	$0$	$0$
$x_2$	$0$	$3/8$	$6/7$	$1$
$x_3$	$1/2$	$1/4$	$1/7$	$0$

.

<sup>3</sup> Another example is bisimulation metrics for probabilistic properties. Because bisimulation is a strong property that can be broken by tiny modifications to transition probabilities, a variety of bisimulation metrics have been proposed, measuring how much bisimulation is violated.

This mapping is called the *abstract channel* denoted by  $\mathbf{C}$ .<sup>4</sup>

It is natural to measure the “vulnerability” of  $X$  using functions  $V_g$  from probability distributions to reals, which are parameterized with a *gain function*  $g : \mathcal{W} \times \mathcal{X} \rightarrow \mathbb{R}$  that models the operational scenario; here  $\mathcal{W}$  is the set of *actions* that the adversary can make, and  $g(w, x)$  is the adversary’s *gain* for doing  $w$  when the secret’s actual value is  $x$ . We then define the  *$g$ -vulnerability* by  $V_g(\pi) = \sup_w \sum_x \pi_x g(w, x)$ , since that is the maximum expected gain over all possible actions. The choice of  $g$  allows us to model a wide variety of operational scenarios; a commonly used one is the *identity* gain function, given by  $\mathcal{W} = \mathcal{X}$  and  $g_{\text{id}}(w, x) = 1$  iff  $w = x$  and 0 otherwise. For this gain function  $V_{g_{\text{id}}}(\pi)$  gives the probability of correctly guessing the secret in one try.

Posterior vulnerability is defined as the average  $g$ -vulnerability in the hyper distribution:  $V_g[\pi \triangleright \mathbf{C}] = \sum_y p(y) V_g(p_{X|y})$ . And then  *$g$ -leakage* is defined as either the *ratio* or the *difference* between the posterior- and prior  $g$ -vulnerability; these are respectively *multiplicative leakage*  $\mathcal{L}_g^\times(\pi, \mathbf{C})$  and *additive leakage*  $\mathcal{L}_g^+(\pi, \mathbf{C})$ .

## 2.2 Metrics

A *proper metric* on a set  $\mathcal{A}$  is a function  $d : \mathcal{A} \rightarrow [0, \infty]$  satisfying the following axioms for all  $x, y, z \in \mathcal{A}$ :

1. **Reflexivity:**  $d(x, x) = 0$ .
2. **Symmetry:**  $d(x, y) = d(y, x)$ .
3. **Anti-symmetry:**  $d(x, y) = d(y, x) = 0 \Rightarrow x = y$ .<sup>5</sup>
4. **Triangle inequality:**  $d(x, z) \leq d(x, y) + d(y, z)$ .
5. **Finiteness:**  $d(x, y) \neq \infty$ .

Various prefixes can be added to “metric” to denote that some of the above conditions are dropped. The following are used in this paper:

1. **Quasi:** symmetry is dropped.
2. **Pseudo:** anti-symmetry is dropped.
3. **Extended:** finiteness is dropped.

Any combination of these prefixes is meaningful. For brevity, we use “metric” to generally refer to *any* of these classes, while we use the exact prefixes (or “proper”) to refer to a precise one.

**Kernel.** The *kernel* of a metric  $d$  is a relation, containing points at distance 0:

$$(x, y) \in \ker(d) \quad \text{iff} \quad d(x, y) = 0.$$

It is easy to see that the metric properties of  $d$  imply the corresponding homonymous properties on relations for  $\ker(d)$  (reflexivity, symmetry, anti-symmetry), while the triangle inequality of  $d$  corresponds to transitivity for  $\ker(d)$ . Hence the kernel of a proper

<sup>4</sup> Because of structural redundancies (e.g. the ordering and labels of columns), distinct channel matrices may denote the same abstract channel.

<sup>5</sup> When  $d$  is not symmetric, this is weaker than the axiom  $d(x, y) = 0 \Rightarrow x = y$  (which is sometimes used in the literature even in the non-symmetric case).

metric is always the *equality relation*  $=$ , the kernel of a pseudometric is an *equivalence relation* (reflexive, symmetric and transitive), and the kernel of a quasimetric is a *partial order* (reflexive, anti-symmetric and transitive).

**Quasimetrics on  $\mathbb{R}, \mathbb{R}_{\geq 0}$ .** On  $\mathbb{R}$  the following quasimetric is of particular interest:

$$\mathbf{q}_{<}^+(x, y) := \max\{y - x, 0\}.$$

Intuitively,  $\mathbf{q}_{<}^+$  measures (additively) “how much smaller” than  $y$  is  $x$ ; 0 means that  $x$  is “no smaller” than  $y$ . Or we can view  $\mathbf{q}_{<}^+$  as measuring “how much  $x \geq y$  is violated”; 0 means that  $x \geq y$  holds (is not violated at all).

On  $\mathbb{R}_{\geq 0}$  we can define a multiplicative variant of  $\mathbf{q}_{<}^+$  as follows:<sup>6</sup>

$$\mathbf{q}_{<}^{\times}(x, y) := \max\{\log_2 \frac{y}{x}, 0\},$$

with the understanding that  $\mathbf{q}_{<}^{\times}(0, y) = 0$  iff  $y = 0$  and  $\infty$  otherwise ( $\mathbf{q}_{<}^{\times}$  is an extended quasimetric). Again,  $\mathbf{q}_{<}^{\times}$  measures “how much smaller” than  $y$  is  $x$ , but this time multiplicatively. Although  $\mathbf{q}_{<}^+$  and  $\mathbf{q}_{<}^{\times}$  are clearly different, they coincide on their kernel:

$$\ker(\mathbf{q}_{<}^+) = \ker(\mathbf{q}_{<}^{\times}) = \geq.$$

Showing the quasimetric properties of  $\mathbf{q}_{<}^+$  is trivial for all but the triangle inequality; for the latter, assuming  $x \leq z$  (the case  $x > z$  is trivial) we have that

$$\begin{aligned} & \mathbf{q}_{<}^+(x, z) \\ = & z - x && \text{“}x \leq z\text{”} \\ = & z - y + y - x \\ \leq & \max\{z - y, 0\} + \max\{y - x, 0\} \\ = & \mathbf{q}_{<}^+(x, y) + \mathbf{q}_{<}^+(y, z). \end{aligned}$$

Given a function  $f : \mathcal{B} \rightarrow \mathcal{A}$  and metric  $d$  on  $\mathcal{A}$  lets us define a metric  $d \circ f$  on  $\mathcal{B}$  by  $(d \circ f)(x, y) := d(f(x), f(y))$ .

**Proposition 1.** *The composition  $d \circ f$  always preserves all metric properties of  $d$  except anti-symmetry, which is also preserved if  $f$  is injective.*

This construction is useful in that  $\mathbf{q}_{<}^{\times} = \mathbf{q}_{<}^+ \circ \log_2$ , from which it follows that  $\mathbf{q}_{<}^{\times}$  is an extended quasimetric on  $\mathbb{R}_{\geq 0}$ .<sup>7</sup>

### 3 Refinement metrics

To design metrics for QIF, we start from the simple observation that we can *compare vulnerabilities* using any (quasi) metric  $q$  on  $\mathbb{R}_{\geq 0}$ . For instance, *leakage* can be measured by comparing prior and posterior vulnerabilities as follows:

$$\mathcal{L}_g^q(\pi, C) := q(V_g(\pi), V_g[\pi \triangleright C]).$$

<sup>6</sup> The use of logarithm is crucial for reflexivity and the triangle inequality. The choice of the base is arbitrary, but  $\log_2$  is interesting due to the connection with min-entropy.

<sup>7</sup> Technically, to establish this we need to view  $\mathbf{q}_{<}^+$  as an extended quasimetric on  $\mathbb{R} \cup \{-\infty\}$  and  $\log_2$  as an injective function  $\mathbb{R}_{\geq 0} \rightarrow \mathbb{R} \cup \{-\infty\}$ .

The usual *additive* and *multiplicative* variants of leakage can be obtained by instantiating  $q$  by  $\mathbf{q}_{<}^+$ ,  $\mathbf{q}_{<}^\times$  (except that for the latter we get  $\log_2$  of the multiplicative leakage<sup>8</sup>):

$$\mathcal{L}_g^{\mathbf{q}_{<}^+}(\pi, C) = \mathcal{L}_g^+(\pi, \mathbf{C}), \quad \mathcal{L}_g^{\mathbf{q}_{<}^\times}(\pi, C) = \log_2 \mathcal{L}_g^\times(\pi, \mathbf{C}).$$

The conceptual advantage of explicitly using a quasimetric (instead of just the difference or ratio) is twofold: first, we can exploit the metric properties of  $q$  and second, we can treat both leakage variants together.

Continuing this line of reasoning, we can obtain a *metric on channels*, by using  $q$  to compare their posterior vulnerabilities. This brings us to our refinement metric.

**Definition 1.** Given classes  $\mathcal{D} \subseteq \mathbb{D}\mathcal{X}$  of distributions and  $\mathcal{G} \subseteq \mathbb{G}\mathcal{X}$  of gain functions, and metric  $q$  on  $\mathbb{R}_{\geq 0}$ , the refinement metric is defined as

$$\text{ref}_{\mathcal{D}, \mathcal{G}}^q(\mathbf{A}, \mathbf{B}) := \sup_{\pi: \mathcal{D}, g: \mathcal{G}} q(V_g[\pi \triangleright \mathbf{A}], V_g[\pi \triangleright \mathbf{B}]).$$

We write  $\text{ref}_{\pi, \mathcal{G}}^q$  instead of  $\text{ref}_{\{\pi\}, \mathcal{G}}^q$  when  $\mathcal{D} = \{\pi\}$ , and similarly for  $\mathcal{G}$ .

Any metric could be meaningful for  $q$  (for instance symmetric ones). But here we restrict to quasimetrics such that  $\ker(q) = (\geq)$ , which we call *order-respecting*. The standard choices of interest are  $\mathbf{q}_{<}^+$  and  $\mathbf{q}_{<}^\times$ , giving additive and multiplicative comparison, for which we write the refinement metrics as  $\text{ref}_{\mathcal{D}, \mathcal{G}}^+$  and  $\text{ref}_{\mathcal{D}, \mathcal{G}}^\times$  respectively. (Note however that many results are independent of the choice of  $q$ .)

When  $\mathcal{D}$  contains at least one point prior, or when  $\mathcal{G}$  contains at least one constant  $g$ , then the max in the definition of  $\mathbf{q}_{<}^+$ ,  $\mathbf{q}_{<}^\times$  can be eliminated, giving

$$\begin{aligned} \text{ref}_{\mathcal{D}, \mathcal{G}}^+(\mathbf{A}, \mathbf{B}) &= \sup_{\pi: \mathcal{D}, g: \mathcal{G}} V_g[\pi \triangleright \mathbf{B}] - V_g[\pi \triangleright \mathbf{A}] && \text{and} \\ \text{ref}_{\mathcal{D}, \mathcal{G}}^\times(\mathbf{A}, \mathbf{B}) &= \sup_{\pi: \mathcal{D}, g: \mathcal{G}} \log_2 \frac{V_g[\pi \triangleright \mathbf{B}]}{V_g[\pi \triangleright \mathbf{A}]} . \end{aligned}$$

Intuitively,  $\text{ref}_{\mathcal{D}, \mathcal{G}}^+$ ,  $\text{ref}_{\mathcal{D}, \mathcal{G}}^\times$  measure (additively or multiplicatively) *how robust* it is to replace  $\mathbf{A}$  by  $\mathbf{B}$ , with respect to the vulnerability of the system. A value of 0 means no risk — $\mathbf{B}$  is never worse than  $\mathbf{A}$  (for  $\mathcal{D}, \mathcal{G}$ )— while a positive value means that there is risk, but it might be small. Note that refinement is *asymmetric* by design: replacing  $\mathbf{A}$  by  $\mathbf{B}$  is inherently different than replacing  $\mathbf{B}$  by  $\mathbf{A}$ .

For fixed  $\pi$  and  $g$ , we can obtain  $\text{ref}_{\pi, g}^q$  as the composition

$$\text{ref}_{\pi, g}^q = q \circ V_g[\pi \triangleright \cdot], \tag{1}$$

and then obtain  $\text{ref}_{\mathcal{D}, \mathcal{G}}^q$  as the sup of a family of quasimetrics

$$\text{ref}_{\mathcal{D}, \mathcal{G}}^q = \sup_{\pi: \mathcal{D}, g: \mathcal{G}} \text{ref}_{\pi, g}^q . \tag{2}$$

This brings us to the following proposition.

<sup>8</sup> This is reminiscent of min-entropy; conceptually, however, we did not use  $\log_2$  to convert vulnerabilities to entropies, we just used  $\mathbf{q}_{<}^\times$  to compare vulnerabilities via a metric of multiplicative nature (the log conveniently turns ratios into a metric).

**Proposition 2.**  $\text{ref}_{\mathcal{D},\mathcal{G}}^q$  always inherits reflexivity, symmetry and the triangle inequality from  $q$ , but not necessarily anti-symmetry and finiteness.

*Proof.* From (1) and Prop. 1 we get that  $\text{ref}_{\pi,g}^q$  inherits all metric properties of  $q$  except anti-symmetry ( $V_g[\pi \triangleright \cdot]$  is not injective). Then, from (2), we only need to show that the sup of metrics preserves reflexivity, symmetry and the triangle inequality. The first two are trivial; for the triangle inequality, let  $d = \sup_i d_i$ , we have:

$$\begin{aligned}
& d(x, z) \\
= & \sup_i d_i(x, z) \\
\leq & \sup_i (d_i(x, y) + d_i(y, z)) && \text{“triangle ineq. for } d_i\text{”} \\
\leq & \sup_i d_i(x, y) + \sup_i d_i(y, z) && \text{“}\sup_i (a_i + b_i) \leq \sup_i a_i + \sup_i b_i\text{”} \\
= & d(x, y) + d(y, z) .
\end{aligned}$$

Clearly, anti-symmetry is not always preserved (take  $\mathcal{G}$  to contain only constant gain functions), nor is finiteness (for  $\mathcal{G} = \mathbb{G}\mathcal{X}$  the sup can be  $\infty$  even if  $q$  itself is finite).  $\square$

Hence  $\text{ref}_{\mathcal{D},\mathcal{G}}^+$  and  $\text{ref}_{\mathcal{D},\mathcal{G}}^\times$  always satisfy reflexivity and the triangle inequality (i.e. they are extended quasi pseudo metrics); they might satisfy additional properties for specific choices of  $\mathcal{D}$  and  $\mathcal{G}$ .

### 3.1 Recovering the refinement relation

The following result states that if  $\mathcal{D}$  and  $\mathcal{G}$  are “sufficiently rich”, then we can obtain  $\sqsubseteq$  as the kernel of the corresponding refinement metric  $\text{ref}_{\mathcal{D},\mathcal{G}}^q$ .

**Theorem 1.** Assume that  $q$  is order-respecting,  $\mathcal{D}$  contains a full-support prior, and  $\mathcal{G}$  contains the non-negative, 1-bounded gain functions with finitely many actions. Then

$$\ker(\text{ref}_{\mathcal{D},\mathcal{G}}^q) = \sqsubseteq .$$

*Proof.* (Sketch.) If  $\mathbf{A} \not\sqsubseteq \mathbf{B}$ , then the proof of the Coriaceous Theorem [10, Thm. 9] shows that for the uniform prior  $\pi^u$ , there is a gain function  $g$ , restricted as in the statement of the theorem, such that  $V_g[\pi^u \triangleright \mathbf{A}] < V_g[\pi^u \triangleright \mathbf{B}]$ ; hence if  $q$  is order-respecting we have  $\text{ref}_{\pi^u,\mathcal{G}}^q(\mathbf{A}, \mathbf{B}) > 0$ . And by embedding the prior into  $g$ , we can replace  $\pi^u$  with any full-support prior.  $\square$

Hence the anti-symmetry of  $\sqsubseteq$  on abstract channels is transferred to  $\text{ref}_{\mathcal{D},\mathcal{G}}^q$ :

**Corollary 1.** Under the conditions of Thm. 1,  $\text{ref}_{\mathcal{D},\mathcal{G}}^q$  is anti-symmetric on abstract channels.

### 3.2 Recovering capacity

Since the noninterfering channel  $\mathbb{1}$  always produces a point hyper, capacity is simply given by the corresponding refinement metric between  $\mathbb{1}$  and  $\mathbf{C}$ .

**Proposition 3.** For any  $q$  and  $\mathcal{D}, \mathcal{G}$  we have that  $\mathcal{ML}_{\mathcal{G}}^q(\mathcal{D}, \mathbf{C}) = \text{ref}_{\mathcal{D},\mathcal{G}}^q(\mathbb{1}, \mathbf{C})$ .

*Proof.* Direct consequence of the fact that  $V_g[\pi \triangleright \mathbb{1}] = V_g(\pi)$ .  $\square$

Recall that using the quasimetrics  $\text{ref}_{\mathcal{D},\mathcal{G}}^+$  or  $\text{ref}_{\mathcal{D},\mathcal{G}}^\times$  for  $q$  make  $\text{ref}_{\mathcal{D},\mathcal{G}}^q$  itself non-symmetric. In particular  $\text{ref}_{\mathcal{D},\mathcal{G}}^q(\mathbb{1}, \mathbf{C})$  gives the capacity of  $\mathbf{C}$ , while  $\text{ref}_{\mathcal{D},\mathcal{G}}^q(\mathbf{C}, \mathbb{1})$  is always zero since  $\mathbf{C} \sqsubseteq \mathbb{1}$ .

A corollary of the above proposition is that the refinement metric of  $\mathbf{A}, \mathbf{B}$  can never exceed the capacity of  $\mathbf{B}$ . This is intuitive since  $\text{ref}_{\mathcal{D},\mathcal{G}}^q(\mathbf{A}, \mathbf{B})$  measures “how much worse than  $\mathbf{A}$  can  $\mathbf{B}$  be”, while the capacity measures “how much worse than  $\mathbb{1}$  (the best possible channel) can  $\mathbf{B}$  be”.

**Corollary 2.** *For any order-respecting  $q$  and  $\mathcal{D}, \mathcal{G}, \mathbf{A}, \mathbf{B}$  we have that*

$$\text{ref}_{\mathcal{D},\mathcal{G}}^q(\mathbf{A}, \mathbf{B}) \leq \mathcal{ML}_{\mathcal{G}}^q(\mathcal{D}, \mathbf{B}).$$

*Proof.* We have that:

$$\begin{aligned} & \text{ref}_{\mathcal{D},\mathcal{G}}^q(\mathbf{A}, \mathbf{B}) \\ \leq & \text{ref}_{\mathcal{D},\mathcal{G}}^q(\mathbf{A}, \mathbb{1}) + \text{ref}_{\mathcal{D},\mathcal{G}}^q(\mathbb{1}, \mathbf{B}) && \text{“triangle ineq. for } \text{ref}_{\mathcal{D},\mathcal{G}}^q\text{”} \\ \leq & \text{ref}_{\mathcal{D},\mathcal{G}}^q(\mathbb{1}, \mathbf{B}) && \text{“} \text{ref}_{\mathcal{D},\mathcal{G}}^q(\mathbf{A}, \mathbb{1}) = 0 \text{ for order-respecting } q\text{”} \\ = & \mathcal{ML}_{\mathcal{G}}^q(\mathcal{D}, \mathbf{B}) && \text{“Prop. 3”} \end{aligned}$$

$\square$

The above corollary in turn implies that  $\text{ref}_{\mathcal{D},\mathcal{G}}^q$  is finite whenever the capacity itself is finite. In particular  $\text{ref}_{\mathbb{D},\mathbb{G}^\dagger}^+$  and  $\text{ref}_{\mathbb{D},\mathbb{G}^+}^\times$  are both finite.

## 4 Computing the additive refinement metric

Given channel matrices  $\mathbf{A}: \mathcal{X} \rightarrow \mathcal{Y}$  and  $\mathbf{B}: \mathcal{X} \rightarrow \mathcal{Z}$  where  $\mathbf{A} \not\sqsubseteq \mathbf{B}$ , we know that sometimes  $\mathbf{B}$  leaks more than  $\mathbf{A}$ . To determine how much worse  $\mathbf{B}$  can be, we wish to compute some version of the refinement metric: here we focus on the additive refinement metric  $\text{ref}_{\pi, \mathbb{G}^\dagger}^+(\mathbf{A}, \mathbf{B})$ , which measures the maximum amount by which the additive leakage of  $\mathbf{B}$  can exceed that of  $\mathbf{A}$  on prior  $\pi$  and over gain functions in  $\mathbb{G}^\dagger \mathcal{X}$ , the class of 1-bounded gain functions.<sup>9</sup>

It is known that we can do this computation via linear programming [4], where the elements of  $\mathbb{G}$  (the matrix representation of  $g$ ) are variables. But here we refine the previously-used approach, enabling it to be done far more efficiently than was previously known. First, to ensure that  $\mathbb{G}$  is in  $\mathbb{G}^\dagger \mathcal{X}$ , we just need to constrain its entries to be at most 1, and (to ensure that its gain values are non-negative) include a special action  $\perp$  whose gain values are all 0. A second issue is that the number of possible actions of  $\mathbb{G}$  is not fixed. But it suffices for  $\mathbb{G}$  to have  $|\mathcal{Y}| + |\mathcal{Z}| + 1$  actions, since this allows

<sup>9</sup> These are the gain functions whose gain values are at most 1 and which always produce non-negative vulnerabilities (even though some of their gain values can be negative). We need to restrict to bounded gain functions to get interesting results, since otherwise we could get arbitrarily large leakage differences simply by scaling up the gain values, making the additive refinement metric infinite.



a different action for each output of  $\mathbf{A}$  and  $\mathbf{B}$ , along with the special  $\perp$  action. Now we recall the trace-based formulation of posterior vulnerability:

$$V_g[\pi \triangleright \mathbf{A}] = \max_{\mathbf{S}^{\mathbf{A}}} \text{tr}(\mathbf{G}\mathbf{D}^\pi \mathbf{A}\mathbf{S}^{\mathbf{A}}) \quad \text{and} \quad V_g[\pi \triangleright \mathbf{B}] = \max_{\mathbf{S}^{\mathbf{B}}} \text{tr}(\mathbf{G}\mathbf{D}^\pi \mathbf{B}\mathbf{S}^{\mathbf{B}})$$

where  $\mathbf{S}^{\mathbf{A}}: \mathcal{Y} \rightarrow \mathcal{W}$  and  $\mathbf{S}^{\mathbf{B}}: \mathcal{Z} \rightarrow \mathcal{W}$  are strategies mapping outputs of  $\mathbf{A}$  and  $\mathbf{B}$  to actions. While the previously-used approach was to try exponentially-many strategies  $\mathbf{S}^{\mathbf{A}}$  and  $\mathbf{S}^{\mathbf{B}}$ , we now argue that this is actually unnecessary. For we can assume without loss of generality that the first  $|\mathcal{Y}|$  rows of  $\mathbf{G}$  contain, in order, optimal actions for the columns of  $\mathbf{A}$ ; the next  $|\mathcal{Z}|$  rows of  $\mathbf{G}$  contain, in order, optimal actions for the columns of  $\mathbf{B}$ ; and finally the last row of  $\mathbf{G}$  is all 0, for the  $\perp$  action. Achieving this just requires reordering the rows of  $\mathbf{G}$  and possibly duplicating some actions (since the same action might be optimal for more than one column). Once this is done, we can use a fixed  $\mathbf{S}^{\mathbf{A}}$  that maps column  $j$  of  $\mathbf{A}$  to row  $j$  of  $\mathbf{G}$ , and a fixed  $\mathbf{S}^{\mathbf{B}}$  that maps column  $k$  of  $\mathbf{B}$  to row  $k+|\mathcal{Y}|$  of  $\mathbf{G}$ . Now we can solve a *single* linear programming problem for that  $\mathbf{S}^{\mathbf{A}}$  and  $\mathbf{S}^{\mathbf{B}}$ :<sup>10</sup>

Choose  $\mathbf{G}$  to maximize  $\text{tr}(\mathbf{G}\mathbf{D}^\pi \mathbf{B}\mathbf{S}^{\mathbf{B}}) - \text{tr}(\mathbf{G}\mathbf{D}^\pi \mathbf{A}\mathbf{S}^{\mathbf{A}})$  subject to  $\mathbf{G}$  having elements of at most 1 and a final all-zero row, and  $\text{opt}(\mathbf{S}^{\mathbf{A}})$ ,

where  $\text{opt}(\mathbf{S}^{\mathbf{A}})$  constrains  $\mathbf{S}^{\mathbf{A}}$  to be optimal for  $\mathbf{G}$ . Remarkably, this means that a maximizing gain function  $\mathbf{G}$  can be found in *polynomial time*.

Let us consider some examples. Consider the following channels:

$\mathbf{A}$	$y_1$	$y_2$	$y_3$
$x_1$	0.1	0.2	0.7
$x_2$	0.3	0.5	0.2
$x_3$	0.6	0.1	0.3

$\mathbf{B}$	$y_1$	$y_2$	$y_3$
$x_1$	0.101	0.200	0.699
$x_2$	0.298	0.501	0.201
$x_3$	0.600	0.099	0.301

where  $\mathbf{B}$  is formed by slightly tweaking the probabilities in  $\mathbf{A}$ . First, given  $\pi$  we can compute the additive capacity of  $\mathbf{A}$ ,  $\mathcal{M}_{\mathbb{G}^\dagger}^+(\pi, \mathbf{A}) = \text{ref}_{\pi, \mathbb{G}^\dagger}^+(\mathbb{1}, \mathbf{A})$ . As predicted in [5, Thm. 6], we find that for any full-support  $\pi$ , the capacity is 0.6, which is indeed 1 minus the sum of the column minimums of  $\mathbf{A}$ ; moreover, the gain function that realizes the capacity is indeed the complement of the  $\pi$ -reciprocal gain function. For instance, when  $\pi = (1/3, 1/2, 1/6)$  it is

$\mathbf{G}$	$x_1$	$x_2$	$x_3$
$w_1$	-2	1	1
$w_2$	1	-1	1
$w_3$	1	1	-5
$\perp$	0	0	0

Next, because  $\mathbf{B}$  is so close to  $\mathbf{A}$ , we expect that its additive refinement metric should be small. On the uniform prior  $\pi^u = (1/3, 1/3, 1/3)$ , we find indeed that  $\text{ref}_{\pi^u, \mathbb{G}^\dagger}^+(\mathbf{A}, \mathbf{B}) =$

<sup>10</sup> Note that deciding whether  $A \sqsubseteq B$  can also be done using a single linear program; but here we achieve more: when  $A \not\sqsubseteq B$  we know that  $B$  can leak more than  $A$ , but we want to also compute *how much* more.

$1063/255000 \approx 0.00416863$ , realized on the gain function

G	$x_1$	$x_2$	$x_3$
$w_1$	$72/85$	$-88/255$	$8/255$
$w_2$	1	$-7/3$	1
$w_3$	1	1	-7
$\perp$	0	0	0

As another example, consider the following channels from [4]:

A	$y_1$	$y_2$	$y_3$	$y_4$
$x_1$	0.1	0.4	0.1	0.4
$x_2$	0.2	0.2	0.3	0.3
$x_3$	0.5	0.1	0.1	0.3

B	$z_1$	$z_2$	$z_3$
$x_1$	0.2	0.22	0.58
$x_2$	0.2	0.4	0.4
$x_3$	0.35	0.4	0.25

They are interesting in that  $A \not\sqsubseteq B$ , but it was historically difficult to find a gain function that makes B leak more than A. On a uniform prior we find that  $\text{ref}_{\pi^u, \mathbb{G}\dagger}^+(\mathbf{A}, \mathbf{B}) = 103/13320 \approx 0.00773273$ , realized on the gain function

G	$x_1$	$x_2$	$x_3$
$w_1$	$-107/148$	$409/444$	$19/444$
$w_2$	$-1/4$	1	$-2/3$
$w_3$	1	$-7/3$	1
$\perp$	0	0	0

To achieve robustness wrt any prior, we might prefer not to compute  $\text{ref}_{\pi, \mathbb{G}\dagger}^+(\mathbf{A}, \mathbf{B})$  for some specific  $\pi$ , but instead to compute  $\text{ref}_{\mathbb{D}, \mathbb{G}\dagger}^+(\mathbf{A}, \mathbf{B})$ . But it is not clear how to do this.

To get some partial insight, we can try an exhaustive search over all priors whose probabilities have a denominator of 1000. We find that the best prior in that set is  $\pi = (624/1000, 29/1000, 347/1000)$ , which gives  $\text{ref}_{\pi, \mathbb{G}\dagger}^+(\mathbf{A}, \mathbf{B}) = 15553/1850000 \approx 0.00840703$ , realized on gain function

G	$x_1$	$x_2$	$x_3$
$w_1$	$-5/7696$	$16/1073$	$4/12839$
$w_2$	$57/208$	1	$-257/347$
$w_3$	1	$-1179/29$	1
$\perp$	0	0	0

Hence we see that unlike the situation of additive capacity where  $\mathcal{ML}_{\mathbb{G}\dagger}^+(\pi, \mathbf{A})$  is the same for any full support  $\pi$ , the additive refinement metric *does* depend on the prior.<sup>11</sup>

<sup>11</sup> There is an interesting phenomenon that we have observed. Given prior  $\pi$ , we can first compute the gain function  $g$  that realizes  $\text{ref}_{\pi, \mathbb{G}\dagger}^+(\mathbf{A}, \mathbf{B})$ , and then we can compute the prior that realizes  $\text{ref}_{\mathbb{D}, g}^+(\mathbf{A}, \mathbf{B})$ . In every case that we have tried, the prior that is found is exactly the  $\pi$  that we started with, suggesting that  $\pi$  is somehow “encoded” into the gain function  $g$ .

## 5 Bounding the additive refinement metric

In the previous section we showed that computing the additive refinement metric for a fixed prior can be done by solving a *single* linear program. Although this gives us a solution in time polynomial in the size of  $\mathbb{C}$ , this solution is still practically feasible only for very modest sizes. In this section, we study efficient techniques for bounding the additive metric based on the Earth Mover’s distance, inspired by the use of the same technique for computing capacity [2].

Moreover, we discuss a simple bound on the additive refinement metric when maximizing over both priors and gain functions. The existence of an efficient exact algorithm for this case remains unknown.

### 5.1 The Earth Mover’s distance

The *Earth Mover’s distance* (EMD), also known as *Wasserstein distance*, is a fundamental distance between *probability distributions*, with numerous applications in computer science. EMD gives a metric on the set  $\mathbb{D}\mathcal{A}$  of probability distributions over an underlying set  $\mathcal{A}$ , equipped with its own “*ground*” metric  $d$ . As a consequence, it can be thought of as a mapping  $\mathbb{W} : \mathbb{M}\mathcal{A} \rightarrow \mathbb{M}\mathbb{D}\mathcal{A}$ , lifting a metric on  $\mathcal{A}$  to a metric on  $\mathbb{D}\mathcal{A}$ .

The idea behind EMD is that the distance between two distributions  $\alpha, \alpha' : \mathbb{D}\mathcal{A}$  is the minimum *cost* of *transforming*  $\alpha$  into  $\alpha'$ . The transformation is performed by moving probability mass between elements  $a, a' : \mathcal{A}$ , where  $\alpha_a$  is seen as the probability mass *available* at  $a$  (supply), and  $\alpha'_{a'}$  as the probability mass *needed* at  $a'$  (demand). The ground metric  $d$  determines the *cost* of this transportation: moving a unit of probability from  $a$  to  $a'$  costs  $d(a, a')$ . There are many strategies for achieving this transformation, the cost of the best one gives the distance between  $\alpha$  and  $\alpha'$ .

More precisely, an earth-moving *strategy* is a joint distribution  $S \in \mathbb{D}\mathcal{A}^2$  whose two marginals are  $\alpha$  and  $\alpha'$ , i.e.  $\sum_{a'} S(a, a') = \alpha_a$  and  $\sum_a S(a, a') = \alpha'_{a'}$ . Intuitively,  $S_{a,a'}$  is the probability mass moved from  $a$  to  $a'$ , and the requirements on the marginals capture the fact that both supply and demand should be respected. We write  $\mathcal{S}_{\alpha, \alpha'}$  for the set of such strategies. Moreover, we write  $\mathcal{E}_S d := \sum_{a, a'} S_{a, a'} d(a, a')$  for the expected value of  $d$  wrt the distribution  $S$ , in other words the total transportation cost of the strategy  $S$ .

**Definition 2.** *The Earth mover’s distance is the mapping  $\mathbb{W} : \mathbb{M}\mathcal{A} \rightarrow \mathbb{M}\mathbb{D}\mathcal{A}$  given by:*

$$\mathbb{W}(d)(\alpha, \alpha') := \inf_{S: \mathcal{S}_{\alpha, \alpha'}} \mathcal{E}_S d.$$

Note that the properties of  $\mathbb{W}(d)$  directly depend on those of  $d$ ; in particular  $\mathbb{W}(d)$  might not be symmetric unless  $d$  itself is symmetric.

The EMD can be computed via a linear program with  $S$  as variables, since  $\mathcal{E}_S d$  is linear on  $S$  and respecting supply and demand can be expressed by linear constraints. However, more efficient network flow algorithms also exist (e.g. [11]). Note, however, that some of these algorithms require  $d$  to be a proper metric.

## 5.2 Bounding the refinement metric using the EMD

We now turn our attention to the problem of bounding the additive refinement metric between two channels  $\mathbf{A}$  and  $\mathbf{B}$ . To do so, we exploit the view of channels as producing *hyper-distributions*. Recall from §2 that the effect of a channel  $\mathbf{A}$  on the prior  $\pi$  can be viewed as a hyper-distribution  $[\pi \triangleright \mathbf{A}]$ , that is a distribution on the posterior distributions produced by each output of  $\mathbf{A}$ . In other words, hypers have type  $\mathbb{D}\mathbb{D}\mathcal{X}$ .

Since hypers are distributions, we can use the EMD to measure the distance between them. But the EMD requires a ground metric on the underlying space, which in our case is  $\mathcal{A} = \mathbb{D}\mathcal{X}$ . In other words, to measure the distance between hypers we first need to measure the distance between posteriors. It turns out that the appropriate distance for our needs is the following.

**Definition 3.** Let  $\pi, \sigma: \mathbb{D}\mathcal{X}$ . Define the “convex separation” quasimetric  $\mathbf{q}_{\text{cs}}: \mathbb{M}\mathbb{D}\mathcal{X}$  as

$$\mathbf{q}_{\text{cs}}(\sigma, \pi) := \max_{x: [\sigma]} \left( 1 - \frac{\pi_x}{\sigma_x} \right).$$

The usefulness of this quasimetric lies in the fact that it can be shown [5] to express the maximum that  $\sigma, \pi$  can be separated by a 1-bounded gain function, that is

$$\mathbf{q}_{\text{cs}}(\sigma, \pi) = \sup_{g: \mathbb{G}^{\dagger}\mathcal{X}} \mathbf{q}_{<}^+(V_g(\sigma), V_g(\pi)).$$

This in turn means that  $V_g$  is non-expansive wrt  $\mathbf{q}_{\text{cs}}$ , which is important due to the connection between EMD and the Kantorovich metric. The technicalities of this construction are available in the appendix, we only state the most relevant result here.

**Theorem 2.** For any prior  $\pi$  and channels  $\mathbf{A}, \mathbf{B}$  it holds that

$$\text{ref}_{\pi, \mathbb{G}^{\dagger}}^+(\mathbf{A}, \mathbf{B}) \leq \mathbb{W}(\mathbf{q}_{\text{cs}})([\pi \triangleright \mathbf{A}], [\pi \triangleright \mathbf{B}]).$$

Hence, to obtain a bound on the refinement metric, we can compute the hypers  $[\pi \triangleright \mathbf{A}]$  and  $[\pi \triangleright \mathbf{B}]$  (i.e. the corresponding output and posterior distributions) and then employ either linear programming or a network flow algorithm to compute the  $\mathbb{W}(\mathbf{q}_{\text{cs}})$  distance between them. Note that the algorithm for EMD should not require the ground metric to be symmetric, since  $\mathbf{q}_{\text{cs}}$  is not.

## 5.3 More relaxed bounds

As discussed in the previous sections, the EMD is the cost of the *best* transportation strategy to transform one distribution into another. As a consequence, one can clearly obtain a bound by choosing *any* transportation strategy, not necessarily the best. That is, if  $S$  is a strategy for transforming  $[\pi \triangleright \mathbf{A}]$  into  $[\pi \triangleright \mathbf{B}]$ , then the cost (wrt  $\mathbf{q}_{\text{cs}}$ ) of  $S$  provides a bound to the refinement metric between  $\mathbf{A}$  and  $\mathbf{B}$ :

$$\text{ref}_{\pi, \mathcal{G}}^+(\mathbf{A}, \mathbf{B}) \leq \mathcal{E}_S \mathbf{q}_{\text{cs}}.$$

Of course the actual EMD gives the best such bound, but the freedom of choosing *any* strategy can be convenient. In §6 an analytic bound for the Crowds protocol is given by choosing an arbitrary strategy, without having to show that this is the optimal one.

#### 5.4 Maximizing over all priors

In view of the apparent difficulty of computing  $\text{ref}_{\mathbb{D}, \mathbb{G}^\dagger}^+(\mathbf{A}, \mathbf{B})$ , we are led to wonder how much the additive refinement bound can be increased by changing from a uniform prior to a non-uniform prior. We find that the possible increase can be bounded:

**Theorem 3.** *For any channel matrices  $\mathbf{A}: \mathcal{X} \rightarrow \mathcal{Y}$  and  $\mathbf{B}: \mathcal{X} \rightarrow \mathcal{Z}$ , we have*

$$\text{ref}_{\mathbb{D}, \mathbb{G}^\dagger}^+(\mathbf{A}, \mathbf{B}) \leq |\mathcal{X}| \cdot \text{ref}_{\pi^u, \mathbb{G}^\dagger}^+(\mathbf{A}, \mathbf{B}).$$

*Proof.* Let prior  $\pi$  be arbitrary, and suppose that  $\text{ref}_{\pi, \mathbb{G}^\dagger}^+(\mathbf{A}, \mathbf{B})$  is realized on gain function  $g$  with matrix representation  $\mathbf{G}$ , which we assume without loss of generality to include an all-zero row. Next recall the trace-based formulation of posterior  $g$ -vulnerability:  $V_g[\pi \triangleright \mathbf{A}] = \max_{\mathbb{S}} \text{tr}(\mathbf{G}\mathbf{D}^\pi \mathbf{A}\mathbf{S})$ . Observe that we can factor  $\mathbf{D}^\pi = \mathbf{D}^\rho \mathbf{D}^{\pi^u}$ , where vector  $\rho$  is given by  $\rho_x = \pi_x / \pi_x^u = |\mathcal{X}| \cdot \pi_x$ . Now we can translate from  $\pi$  to the uniform prior  $\pi^u$ :

$$V_g[\pi \triangleright \mathbf{A}] = \max_{\mathbb{S}} \text{tr}(\mathbf{G}\mathbf{D}^\pi \mathbf{A}\mathbf{S}) = \max_{\mathbb{S}} \text{tr}(\mathbf{G}\mathbf{D}^\rho \mathbf{D}^{\pi^u} \mathbf{A}\mathbf{S}) = V_{g'}[\pi^u \triangleright \mathbf{A}],$$

where  $g'$  is the gain function with matrix representation  $\mathbf{G}\mathbf{D}^\rho$ . The only problem is that  $g'$  may not be in  $\mathbb{G}^\dagger \mathcal{X}$ . It *does* give non-negative vulnerabilities (since  $\mathbf{G}\mathbf{D}^\rho$  contains an all-zero row), but it may have gain values bigger than 1. We can deal with this by scaling  $g'$  down, dividing its gain values by  $\max_x \pi_x / \pi_x^u = \max_x \pi_x \cdot |\mathcal{X}|$ , which is at most  $|\mathcal{X}|$ . If we let  $g''$  be the scaled-down version of  $g'$ , then we have  $g'' \in \mathbb{G}^\dagger \mathcal{X}$  and  $|\mathcal{X}| \cdot V_{g''}[\pi^u \triangleright \mathbf{A}] = V_g[\pi \triangleright \mathbf{A}]$ . Since the same equality holds for  $\mathbf{B}$ , the desired inequality follows. (Note that it is only an inequality, since  $g''$  may not be optimal for  $\pi^u$ .)  $\square$

Applying this theorem to the last example from §4, where the additive refinement metric on the uniform prior was about 0.0077, we see that the metric over *all* priors is at most 3 times as large, or 0.0232. While this bound is considerably larger than 0.0084 (the maximum found in our partial search), it may still be useful.

## 6 Case study : the Crowds protocol

Crowds is a simple protocol for anonymous web surfing. The context is that a user, called the *initiator*, wants to contact a *web server* but does not want to disclose his identity to the server. That is achieved by collaborating with a group of other users, called the *crowd*, who participate in the protocol to facilitate the task. The protocol is essentially a simple probabilistic routing protocol:

- In the first step, the initiator selects a user uniformly (including possibly himself) and forwards the request to him. The user who receives the message now becomes the (new) *forwarder*.
- A forwarder, upon receiving a message, flips a (biased) probabilistic coin: with probability  $\varphi$  he forwards the message to a new user (again chosen uniformly, including himself), but with remaining probability  $1 - \varphi$  he instead delivers the message directly to the web server.

If  $\varphi < 1$  then with probability 1 the message will eventually arrive at the web server, with the expected number of hops being  $1/(1-\varphi)$ .

Concerning anonymity, from the point of view of the web server all users seem equally likely to have been the initiator (assuming a uniform prior), due to the first step's always being a forwarding step. The analysis becomes more interesting however if we assume that some users in the crowd are corrupt, reporting to an adversary who is trying to discover messages' initiators.<sup>12</sup>

If user  $a$  forwards a message to a corrupted user, we say that  $a$  was *detected*. Since  $a$ 's being detected is more likely when he is the initiator than when he is not, the protocol does leak information: detecting  $a$  creates a posterior where  $a$  has a greater chance of being the initiator than he had *a priori*. Still, if  $\varphi > 0$  then the posterior is *not* a point distribution: user  $a$  can “plead not guilty” by claiming the he was merely forwarding for someone else. The greater  $\varphi$  is, the more plausible that claim becomes: the forwarding probability  $\varphi$  can be seen as trading anonymity for utility (expected number of hops).

In their original work, Reiter and Rubin [12] introduced the property of “probable innocence” to analyze Crowds, meaning, roughly speaking, that each user appears more likely to be innocent than guilty. For our purposes, however, we study Crowds in the context of QIF.<sup>13</sup>

## 6.1 Modeling the Crowds protocol

To perform an analysis of the anonymity guarantees of Crowds, we start by modeling the protocol as an information-theoretic channel. Letting  $n$ ,  $c$ ,  $m$  be the number of honest, corrupted, and total users respectively (i.e.  $n+c = m$ ), the secret and observable events are  $\mathcal{X} := \{x_1, \dots, x_n\}$ ,  $\mathcal{Y} := \{y_1, \dots, y_n, s\}$ , with  $x_a$  denoting that (honest) user  $a$  is the initiator,  $y_b$  denoting that user  $b$  was detected (forwarded a message to a corrupted user), and  $s$  denoting that the message was delivered to the web server without having been forwarded to a corrupted user.

To construct the channel matrix  $\mathbf{C}$ , we need for any honest user  $a$  to compute the probability of producing each observable event, given the secret event  $x_a$ . Due to symmetry, the probability of detecting the initiator, or some non-initiator, is the same for all initiators. In other words, the probabilities  $p(y_b|x_a)$  only depend on whether  $a=b$  or  $a \neq b$ , and not on the exact value of  $a, b$ . Hence the entries of  $\mathbf{C}$  can only have 3 distinct values,  $\mathbf{C}_{x_a, s} = \alpha$ ,  $\mathbf{C}_{x_a, y_a} = \beta$ , and  $\mathbf{C}_{x_a, y_b} = \gamma$ ,  $a \neq b$ , giving the following channel matrix:<sup>14</sup>

$$\mathbf{C} = \begin{matrix} & y_1 & \cdots & y_n & s \\ \begin{matrix} x_1 \\ \vdots \\ x_n \end{matrix} & \begin{bmatrix} \beta & \cdots & \gamma & \alpha \\ \cdot & \ddots & \cdot & \cdot \\ \gamma & \cdots & \beta & \alpha \end{bmatrix} \end{matrix} \quad \text{where} \quad \begin{aligned} \alpha & := \frac{n - \varphi n}{m - \varphi n}, \\ \beta & := \frac{c(m - \varphi(n-1))}{m(m - \varphi n)}, \\ \gamma & := \frac{c\varphi}{m(m - \varphi n)}. \end{aligned}$$

<sup>12</sup> This is a relatively weak adversarial model; a model in which the adversary can view any communication in the network would be too strong for Crowds.

<sup>13</sup> See [3] for a discussion of how probable innocence relates to QIF.

<sup>14</sup> See [3] for the exact derivations.

Finally, given a prior  $\pi$  on  $\mathcal{X}$ , modeling how likely each user is to be the initiator a-priori, we can quantify the anonymity of the protocol by considering its (additive or multiplicative)  $g$ -leakage, for a properly constructed gain function  $g$ . For instance, the identity gain function  $g_{\text{id}}$  is a reasonable choice for this protocol, modeling an adversary that tries to guess the initiator of the message in one try. For this choice of gain function, and the uniform prior  $\pi^u$ , the additive leakage is given by

$$\mathcal{L}_{g_{\text{id}}}^+(\pi^u, \mathbf{C}) = \frac{nc - c}{nc + n} . \quad (3)$$

But  $g_{\text{id}}$  is not the only reasonable choice for this protocol. It could be argued, for instance, that the whole point behind “probable innocence” is that the adversary would not incriminate a user unless there is sufficient evidence — that is, a *wrong* guess should have a *penalty*, compared to *not guessing at all*. Such adversaries can be expressed by suitable gain functions [3].

A crucial question, then, is how can we make *robust* conclusions in our analysis of information leakage, in particular when we tune a parameter of the protocol such as  $\varphi$ ?

## 6.2 Increasing $\varphi$ : refinement relation

A striking fact about Crowds, first observed in [8], is that its Bayes leakage for a uniform prior, given by (3), does not depend on  $\varphi$ ! This is remarkable since the fact of forwarding is the only reason why a detected user can plausibly pretend not be the initiator. Intuitively, the more likely a user is to forward the more likely detecting the wrong user becomes, so anonymity is improved. Understanding this paradox is beyond the scope of this paper (an explanation is given in [3]), but we should note that this phenomenon *only* happens in the very specific case of  $g_{\text{id}}$  with a uniform prior: for most other priors and gain functions, the corresponding leakage does depend on  $\varphi$ , as expected.

This brings us to the natural question: *how* does leakage depend on  $\varphi$ ? Answering this question is of great importance to the designer of the system who tries to optimize the parameters, but without compromising its anonymity.

The easier case to answer is when  $\varphi$  is *increased*. Intuitively, this operation should be *safe*, in the sense that the resulting protocol should leak no more than before. But can we be sure that this *always* happens? Could it be the case that for some strange adversary, modeled by some gain function  $g$ , increasing  $\varphi$  causes the leakage to increase, leading to a less safe protocol?

The theory of refinement provides a robust answer. Let  $\mathbf{C}^\varphi$  denote the channel matrix for a particular choice of  $\varphi$ . It can be shown that

$$\mathbf{C}^\varphi \sqsubseteq \mathbf{C}^{\tilde{\varphi}} \quad \text{iff} \quad \varphi \leq \tilde{\varphi} ,$$

by showing that  $\mathbf{C}^{\tilde{\varphi}} = \mathbf{C}^\varphi \mathbf{R}$  for a properly constructed channel  $\mathbf{R}$ . As a consequence  $\mathbf{C}^{\tilde{\varphi}}$  is a safer protocol; it leaks no more than  $\mathbf{C}^\varphi$  for *all* priors and gain functions.

## 6.3 Decreasing $\varphi$ : refinement metric

Although refinement establishes in a very robust way that increasing  $\varphi$  is a safe operation, it can be argued that the interesting case is exactly the opposite: *decreasing*  $\varphi$ .

This is because a system designer is typically interested in improving the *utility* of the system, which in this case corresponds to its *latency*. Since the expected number of hops to reach the server is  $1/1-\varphi$ , decreasing  $\varphi$  lowers the protocol’s latency, so such a change is desirable.

However, decreasing  $\varphi$  is *not always* safe. Although Bayes leakage for a uniform prior remains unaffected, we can experimentally verify that for all other priors, and for many gain functions, the leakage indeed increases. We thus have a typical trade-off between anonymity and utility; we might be willing to accept a less safe protocol, and *replace  $\tilde{\varphi}$  by a smaller  $\varphi$* , if we knew that the increase in leakage cannot exceed a certain threshold. Can we answer this question in a *robust* way?

**Using the refinement metric.** The refinement metric provides a natural answer to our question. For a fixed  $\pi$ , we know that  $\text{ref}_{\pi, \mathbb{G}^\dagger}^+(\mathbf{C}^{\tilde{\varphi}}, \mathbf{C}^\varphi)$  gives the maximum increase in vulnerability when we replace  $\tilde{\varphi}$  by  $\varphi$ , for any 1-bounded gain function  $g: \mathbb{G}^\dagger$ .

The metric can be computed using the linear programming technique of §4). As an example, for  $n = 110$ ,  $c = 5$ ,  $\varphi = 0.55$  and  $\tilde{\varphi} = 0.6$  we find that the refinement metric is approximately 0.0076. This means that lowering the probability of forwarding from 0.6 to 0.55 cannot have a huge effect on the vulnerability; for instance if the adversary’s probability of a correct guess (Bayes vulnerability) was  $p$ , it can be at most  $p + 0.0076$  after the change. Note that  $\varphi = 0.55$  gives 12.5% *lower latency* compared to  $\varphi = 0.6$ , so we might decide that the benefit outweighs the risk.

If computing the exact metric is not feasible, we might still be able to apply the EMD bound discussed in §5, given by the EMD (wrt  $\mathbf{q}_{\text{cs}}$ ) between the hypers  $[\pi^u \triangleright \mathbf{C}^{\tilde{\varphi}}]$  and  $[\pi^u \triangleright \mathbf{C}^\varphi]$ . In our example, this is approximately 0.017; it is more than twice as large as the real metric, but it still provides a useful bound on the loss of anonymity.

**An analytic bound for Crowds.** As discussed in §5, an upper bound to the refinement metric can be obtained by computing, not the actual EMD, but the cost (wrt  $\mathbf{q}_{\text{cs}}$ ) of *any transportation strategy* from  $[\pi \triangleright \mathbf{C}^{\tilde{\varphi}}]$  to  $[\pi \triangleright \mathbf{C}^\varphi]$ . For the uniform prior,  $[\pi^u \triangleright \mathbf{C}^\varphi]$  can be computed analytically; then we can choose any strategy to obtain an analytic bound.

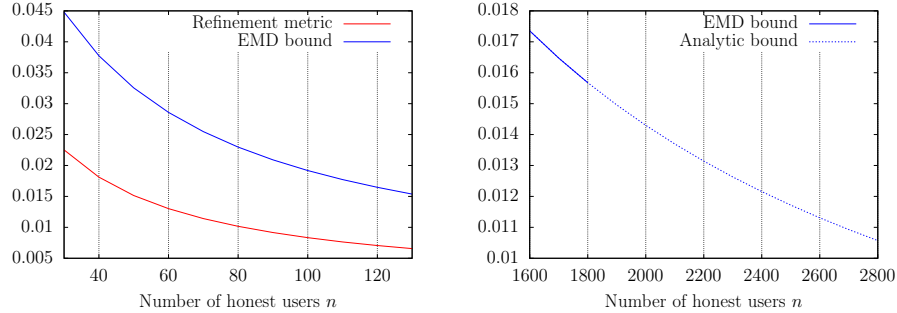
Following the calculations of [3], the hyper distribution  $[\pi^u \triangleright \mathbf{C}^\varphi]$  has

$$\begin{aligned} \text{output distribution} \quad \rho &= (1-\alpha/n, \dots, 1-\alpha/n, \alpha), \\ \text{and posteriors} \quad \delta^{y_b} &= \left( \frac{\varphi}{m}, \dots, \frac{\varphi}{m}, \frac{m-\varphi(n-1)}{m}, \frac{\varphi}{m}, \dots, \frac{\varphi}{m} \right), \\ \delta^s &= (1/n, \dots, 1/n). \end{aligned}$$

Note that the detection events  $y_b$  produce a posterior  $\delta^{y_b}$  assigning higher probability to the detected user  $b$ , while  $s$  provides no information, producing a uniform posterior. Denoting by  $\tilde{\alpha}, \tilde{\rho}, \dots$  the corresponding quantities for  $\tilde{\varphi} \geq \varphi$ , we can check that  $\tilde{\alpha} \leq \alpha$ ; hence  $\tilde{\rho}$  has lower probability than  $\rho$  on its last element and higher on its first  $n$  elements. A reasonable transportation strategy<sup>15</sup>  $S$  from  $[\pi \triangleright \mathbf{C}^{\tilde{\varphi}}]$  to  $[\pi \triangleright \mathbf{C}^\varphi]$  is

<sup>15</sup> Experiments suggest that  $S$  might be optimal, although this is neither required nor proved.





**Fig. 1.** Refinement metric, EMD and analytic bounds for Crowds

then to move all probability mass from  $\tilde{\delta}^s$  to  $\delta^s$ , while splitting the probability of  $\tilde{\delta}^{y_b}$  between  $\delta^{y_b}$  and  $\delta^s$ . This strategy produces the following bound.

**Theorem 4.** For any  $0 < \varphi \leq \tilde{\varphi} \leq 1$  it holds that:

$$\text{ref}_{\pi^u, \mathbb{G}^\dagger}^+(\mathbf{C}^{\tilde{\varphi}}, \mathbf{C}^\varphi) \leq (1 - \alpha)\left(1 - \frac{\varphi}{\tilde{\varphi}}\right) + (\alpha - \alpha') \left(1 - \frac{m}{mn - \tilde{\varphi}n(n-1)}\right).$$

**Experimental evaluation** All techniques presented in this paper have been implemented in the libqif library [1]. An experimental evaluation of the additive refinement metric for Crowds is shown in Fig. 1. On the left-hand side, we compute the exact metric as well as the EMD bound between two instances of Crowds with  $\varphi = 0.55$  and  $\tilde{\varphi} = 0.6$ , for a uniform prior,  $c = 5$  corrupted users, while varying the number of honest users  $n$ . We see that as  $n$  increases  $\text{ref}_{\pi^u, \mathbb{G}^\dagger}^+(\mathbf{C}^{\tilde{\varphi}}, \mathbf{C}^\varphi)$  becomes small, guaranteeing that the decrease of  $\varphi$  has limited effect on the leakage. The EMD bound is not tight, but as  $n$  increases it becomes small enough to provide meaningful guarantees.

Computing the exact metric was possible for up to  $n = 130$  within a computation threshold of 10 minutes. Computing the EMD bound, on the other hand, was feasible for much larger sizes. On the right-hand side of Fig. 1 we show both the EMD and the analytic bound of the previous section, for the same example, but for larger numbers of users and  $c = 70$ . The two bounds coincide, suggesting that the transportation strategy chosen for the analytic bound is likely the optimal one. We were able to compute the EMD bound for up to  $n = 1800$ , while the analytic bound is applicable to any size. Note that, although the exact refinement metric is unknown, the bounds are sufficiently small to provide meaningful guarantees about the anonymity of the protocol.

## References

1. libqif : Quantitative information flow c++ library. <https://github.com/chatziko/libqif>
2. Alvim, M.S., Chatzikokolakis, K., McIver, A., Morgan, C., Palamidessi, C., Smith, G.: Additive and multiplicative notions of leakage, and their capacities. In: Proc. 27th IEEE Computer Security Foundations Symposium (CSF 2014). pp. 308–322 (2014)

3. Alvim, M.S., Chatzikokolakis, K., McIver, A., Morgan, C., Palamidessi, C., Smith, G.: The Science of Quantitative Information Flow. Springer (2019)
4. Alvim, M.S., Chatzikokolakis, K., Palamidessi, C., Smith, G.: Measuring information leakage using generalized gain functions. In: Proc. 25th IEEE Computer Security Foundations Symposium (CSF 2012). pp. 265–279 (Jun 2012)
5. Chatzikokolakis, K.: On the additive capacity problem for quantitative information flow. In: Proc. 15th International Conference on Quantitative Evaluation of Systems (QEST). pp. 1–19 (2018)
6. Clark, D., Hunt, S., Malacaria, P.: Quantitative analysis of the leakage of confidential data. In: Proc. Workshop on Quantitative Aspects of Programming Languages. Electr. Notes Theor. Comput. Sci, vol. 59 (3), pp. 238–251 (2001)
7. Clarkson, M., Myers, A., Schneider, F.: Belief in information flow. In: Proc. 18th IEEE Computer Security Foundations Workshop (CSFW '05). pp. 31–45 (2005)
8. Espinoza, B., Smith, G.: Min-entropy as a resource. Information and Computation (Special Issue on Information Security as a Resource) 226, 57–75 (Apr 2013)
9. Köpf, B., Basin, D.: An information-theoretic model for adaptive side-channel attacks. In: Proc. 14th ACM Conference on Computer and Communications Security (CCS '07). pp. 286–296 (2007)
10. McIver, A., Morgan, C., Smith, G., Espinoza, B., Meinicke, L.: Abstract channels and their robust information-leakage ordering. In: Proc. 3rd Conference on Principles of Security and Trust (POST 2014). pp. 83–102 (2014)
11. Pele, O., Werman, M.: Fast and robust earth mover’s distances. In: ICCV (2009)
12. Reiter, M.K., Rubin, A.D.: Crowds: Anonymity for web transactions. ACM Transactions on Information Systems Security 1(1), 66–92 (1998)
13. Villani, C.: Topics in optimal transportation. No. 58, American Mathematical Soc. (2003)

## A Proofs of Section 5

In this section we provide the proof of the main result of Section 5.

### A.1 The Kantorovich distance

The Kantorovich metric plays a fundamental role in our bounding technique. Given metrics  $d_A : \mathbb{M}\mathcal{A}$ ,  $d_B \in \mathbb{M}\mathcal{B}$ , a function  $f : \mathcal{A} \rightarrow \mathcal{B}$  is called Lipschitz wrt  $d_A, d_B$  iff

$$d_B(f(a), f(a')) \leq d_A(a, a') \quad \text{for all } a, a' \in \mathcal{A}.$$

Let  $\mathbb{C}^{d_A, d_B}\mathcal{A}$  denote the set of all such functions.

The standard Kantorovich construction transforms metrics on some set  $\mathcal{A}$ , to metrics on  $\mathbb{D}\mathcal{X}$ , i.e. metrics measuring the distance between two probability distributions on  $\mathcal{A}$ . Formally, it is the lifting  $\mathbb{K} : \mathbb{M}\mathcal{A} \rightarrow \mathbb{M}\mathbb{D}\mathcal{A}$  given by

$$\mathbb{K}(d)(\alpha, \alpha') := \sup_{F: \mathbb{C}^{d, d_{\mathbb{R}}}\mathcal{A}} |\mathcal{E}_{\alpha}F - \mathcal{E}_{\alpha'}F|.$$

Note the two uses of the Euclidean distance  $d_{\mathbb{R}}$  in the above definition: first, we consider Lipschitz functions wrt to  $d$  (the metric being lifted) and  $d_{\mathbb{R}}$ ; second, we use  $d_{\mathbb{R}}$  to compare the expected values of these functions.

For the purposes of bounding the refinement metric, it is more convenient to use a *quasimetric variant* of Kantorovich, in which  $d_{\mathbb{R}}$  is replaced by  $\mathbf{q}_{<}^+$ .

**Definition 4.** *The Kantorovich quasimetric is the mapping  $\mathbb{K}^+ : \mathbb{M}\mathcal{A} \rightarrow \mathbb{M}\mathbb{D}\mathcal{A}$  given by:*

$$\mathbb{K}^+(d)(\alpha, \alpha') := \sup_{F: \mathbb{C}^{d, \mathfrak{q}_{\leq}^+ \mathcal{A}}} \mathfrak{q}_{\leq}^+(\mathcal{E}_{\alpha} F, \mathcal{E}_{\alpha'} F) = \sup_{F: \mathbb{C}^{d, \mathfrak{q}_{\leq}^+ \mathcal{A}}} \mathcal{E}_{\alpha'} F - \mathcal{E}_{\alpha} F .$$

Note that the  $\max$  in the definition of  $\mathfrak{q}_{\leq}^+$  is removed since the  $\sup$  is anyway non-negative.

An important property of the Kantorovich metric is its *dual* formulation as the *Earth Mover's distance*  $\mathbb{W}(d)$  (Def. 2, also known as Wasserstein metric). The well-known Kantorovich-Rubinstein theorem states that  $\mathbb{K}(d) = \mathbb{W}(d)$  (when  $d$  is a proper metric and satisfies further assumptions, see below). In our case, however, we use  $\mathbb{K}^+$  *instead of*  $\mathbb{K}$ , and moreover we apply it to distances  $d$  that are *not symmetric*. It turns out that the full Kantorovich-Rubinstein theorem still holds in our case (under the same conditions).

**Theorem 5 (Kantorovich-Rubinstein, extended).** *Let  $\mathcal{A}$  be a Polish (i.e. complete and separable) metric space, let  $d$  be a lower semi-continuous distance function on  $\mathcal{A}$  satisfying reflexivity and the triangle inequality. Then  $\mathbb{K}^+(d) = \mathbb{W}(d)$ .*

*Proof.* In [13, Thm. 1.3, page 19], a “Kantorovich duality” is proven, involving generic (lower semi-continuous) cost functions. Then, Kantorovich-Rubinstein is retrieved as a spacial case when the cost function is a metric [13, Thm. 1.14, page 34]. In the proof of the latter (which is quite short), only the reflexivity and the triangle inequality of  $d$  are used. Moreover, the change of  $d_{\mathbb{R}}$  to  $\mathfrak{q}_{\leq}^+$  in the definition of  $\mathbb{K}^+$  does not affect the proof.  $\square$

A surprising consequence of the above result is that, under the assumptions of Kantorovich-Rubinstein, we have:

$$\mathbb{K}^+(d) = \mathbb{K}(d) = \mathbb{W}(d) .$$

Note that the classes of Lipschitz functions in  $\mathbb{K}^+$  and  $\mathbb{K}$  are different, still maximizing over them turns out to give the same result! So, although the formulation of  $\mathbb{K}^+$  is essential for computing capacity (see the next section), in fact it gives the standard Kantorovich construction.

## A.2 Using Kantorovich/EMD to bound the refinement metric

The connection between Kantorovich and EMD allows us to obtain a bound on the refinement metric. To do so, we use the fact that posterior vulnerability, for any 1-bounded gain function, is non-expansive wrt convex separation quasimetric (Def. 3).

**Lemma 1 ([5]).** *For any  $g \in \mathbb{G}^{\dagger} \mathcal{X}$ ,  $V_g$  is Lipschitz wrt  $\mathfrak{q}_{\text{cs}}, d_{\mathbb{R}}$ .*

Since Kantorovich distance is given by maximizing over *any* Lipschitz function, it clearly gives an upper bound on the refinement metric. And the connection between Kantorovich and EMD allows us to efficiently compute this bound.

**Theorem 2.** *For any prior  $\pi$  and channels  $\mathbf{A}, \mathbf{B}$  it holds that*

$$\text{ref}_{\pi, \mathbb{G}^\dagger}^+(\mathbf{A}, \mathbf{B}) \leq \mathbb{W}(\mathbf{q}_{\text{cs}})([\pi \triangleright \mathbf{A}], [\pi \triangleright \mathbf{B}]) .$$

*Proof.* From Lemma 1 we get that  $\text{ref}_{\pi, \mathbb{G}^\dagger}^+(\mathbf{A}, \mathbf{B}) \leq \mathbb{K}^+(\mathbf{q}_{\text{cs}})([\pi \triangleright \mathbf{A}], [\pi \triangleright \mathbf{B}])$ , and the result follows from the (extended) Kantorovich-Rubinstein theorem.