

LSH Ensemble: Internet-Scale Domain Search

Erkang (Eric) Zhu⁺, Fatemah Nargesian⁺, Ken Q. Pu^{*}, Renée J. Miller⁺

+ University of Toronto * University of Ontario Institute of Technology

Alice, a data scientist, wishes to understand the factors affecting companies' funding to the universities.

Scenario 1: she has the two tables at hand.

Industry Partners	Province	Grant Amount
NVIDIA	Ontario	...
Imperial Oil Ltd	Alberta	...
Hydro-Qubec	Quebec	...
...



Solution 1: join the two tables on the Industry Partners and Company columns.

Company	CRA Tax ID	Revenue
NVIDIA	C0112	...
Imperial Oil Ltd	C1234	...
IBM Canada Ltd	C5678	...
...

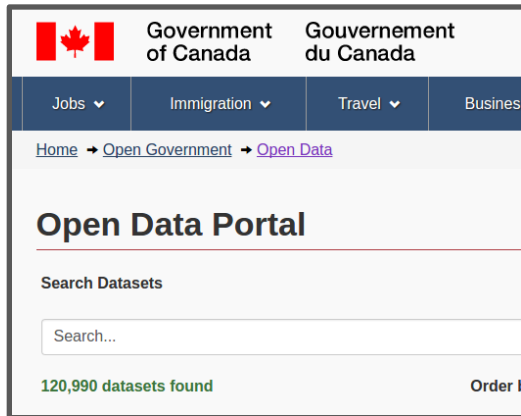
Scenario 2: she has only one table, and a **repository** of MANY tables - manual inspection is not preferred.

Industry Partners	Province	Grant Amount
NVIDIA	Ontario	...
Imperial Oil Ltd	Alberta	...
Hydro-Qubec	Quebec	...
...

Solution 2: requires a search engine for relevant datasets.



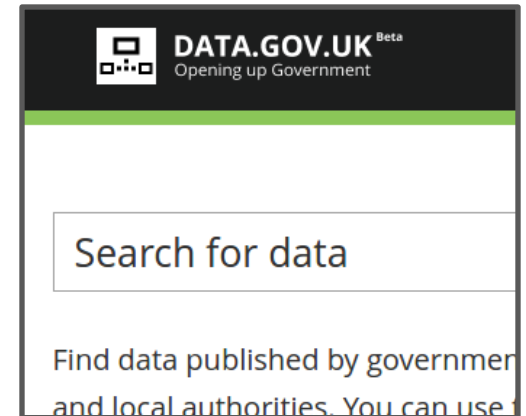
An Internet of Data



The screenshot shows the Government of Canada Open Data Portal. At the top, there is a header with the Canadian flag and the text "Government of Canada" and "Gouvernement du Canada". Below the header is a navigation menu with categories: "Jobs", "Immigration", "Travel", and "Business". A breadcrumb trail reads "Home → Open Government → Open Data". The main heading is "Open Data Portal". Underneath, there is a "Search Datasets" section with a search input field containing "Search...". Below the search field, it states "120,990 datasets found" and "Order by".

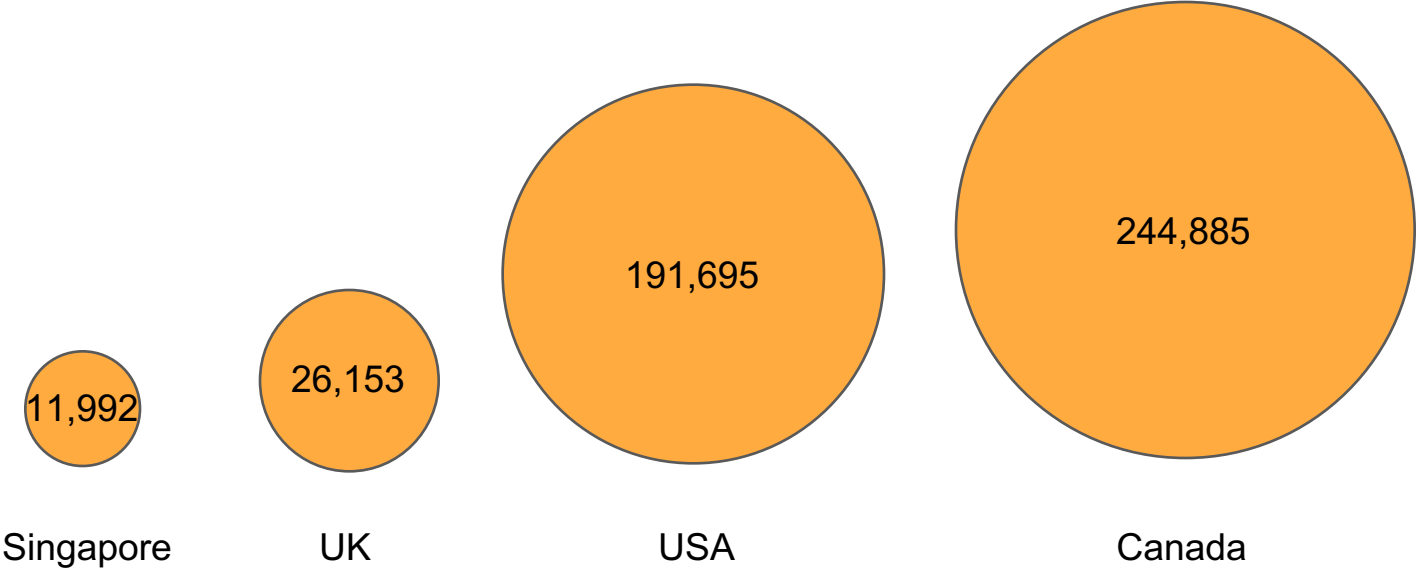


The screenshot shows the U.S. Data.gov website. The header features the "DATA.GOV" logo with the American flag and navigation links for "DATA", "TOPICS", "IMPACT", and "APPLICATIONS". The main heading is "The home of the U.S. Government's open data". Below this, a sub-heading reads "Here you will find data, tools, and resources to conduct research, develop applications, design data visualizations, and [more](#)." A prominent "GET STARTED" button is displayed, with the text "SEARCH OVER 185,973 DATASETS" underneath. A search bar at the bottom contains the text "Credit Card Complaints".



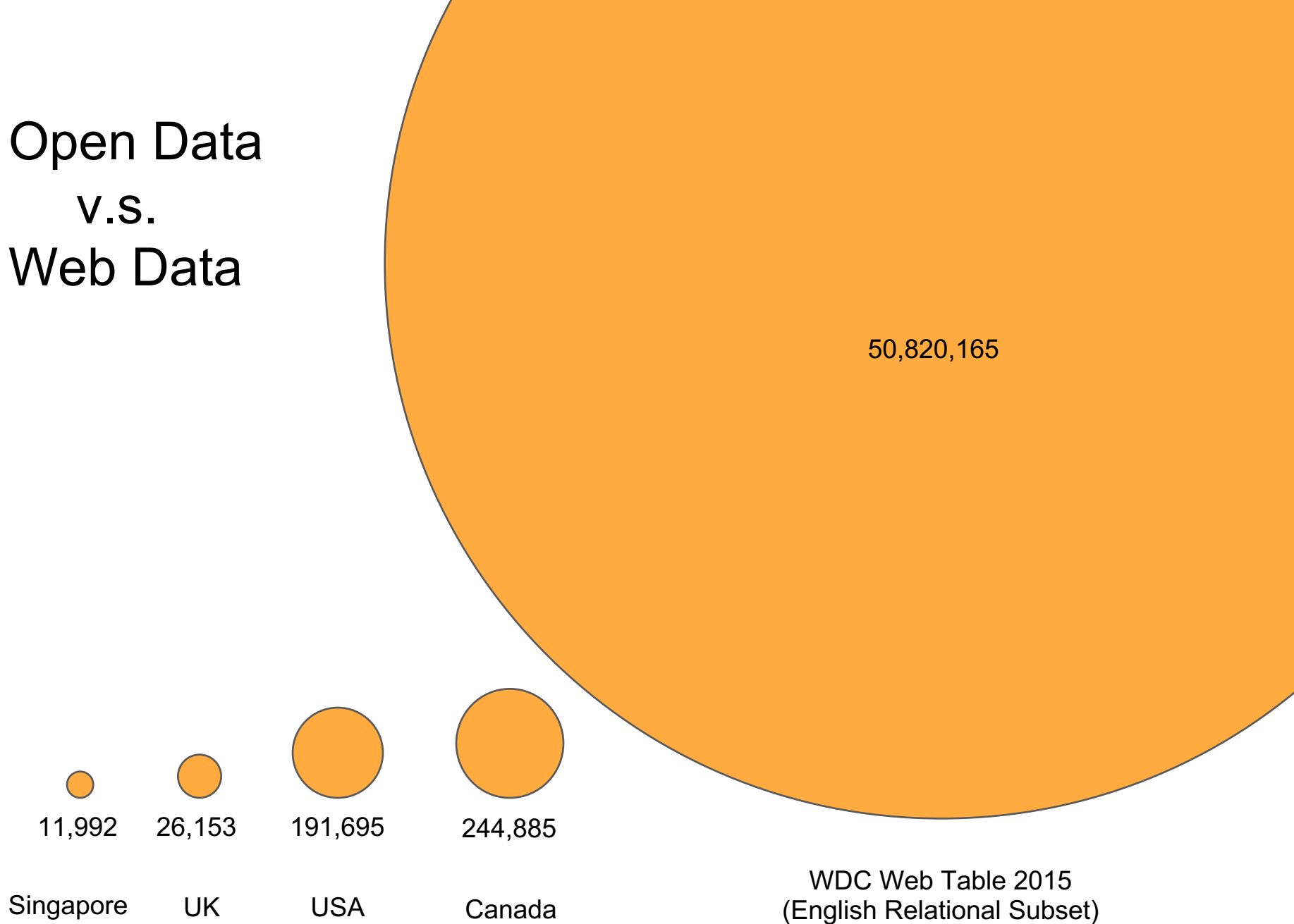
The screenshot shows the U.K. Data.gov.uk website. The header features the "DATA.GOV.UK" logo with the tagline "Opening up Government" and a "Beta" badge. The main heading is "Search for data". Below this, a sub-heading reads "Find data published by government and local authorities. You can use".

Examples of Open Data (Early 2016)

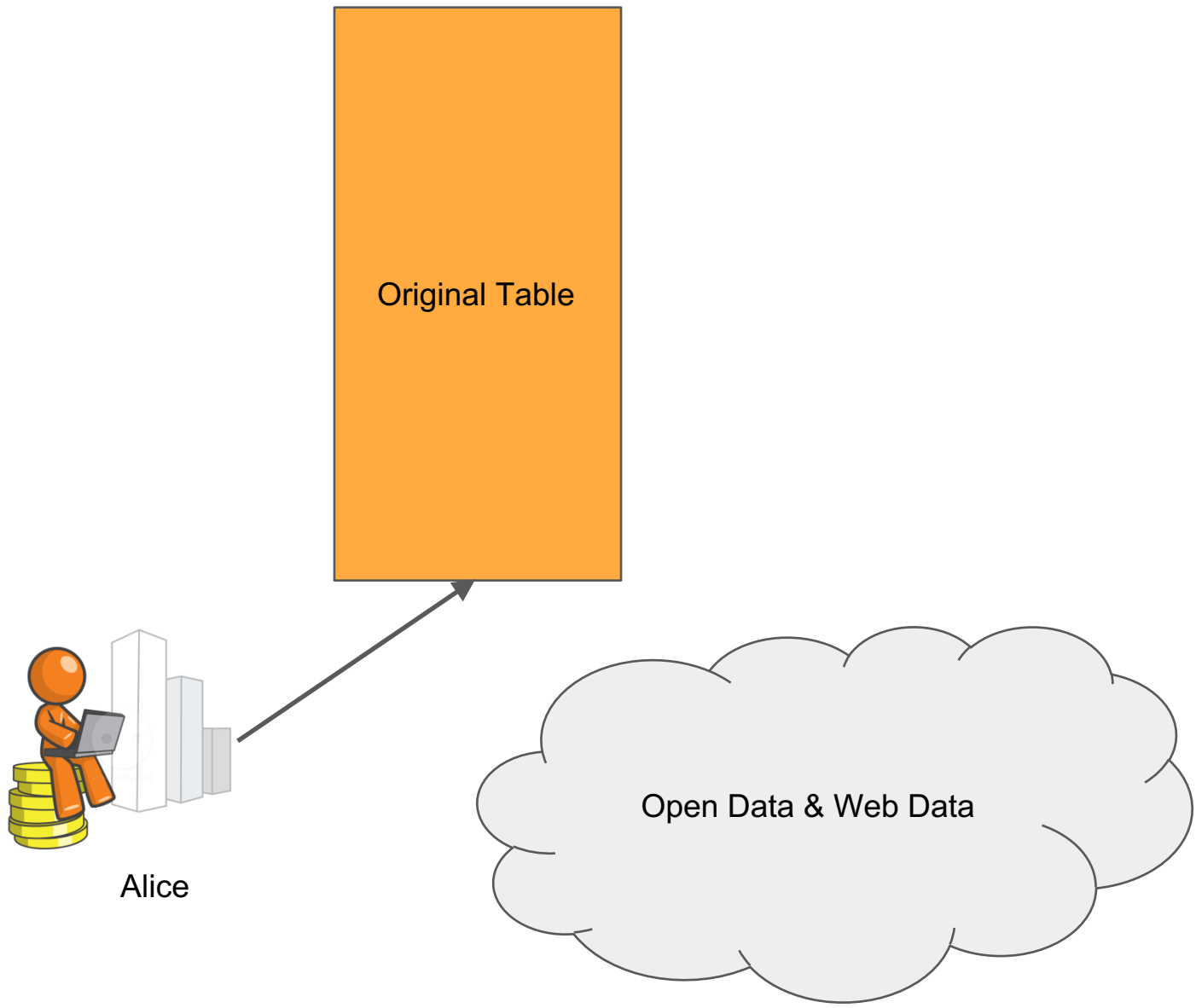


Number of Datasets by Country

Open Data v.s. Web Data



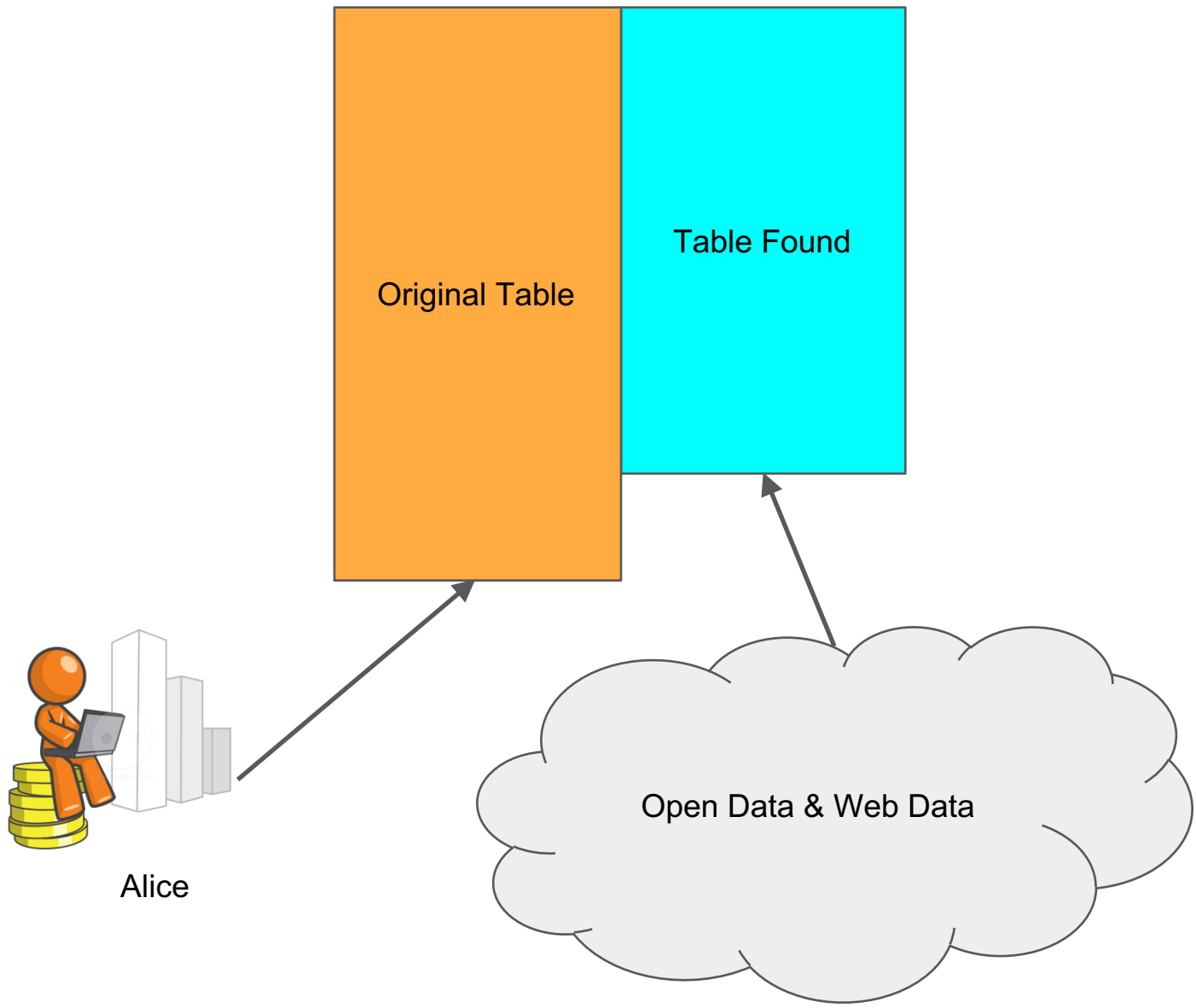
What is a reasonable objective for searching datasets?



Original Table

Alice

Open Data & Web Data

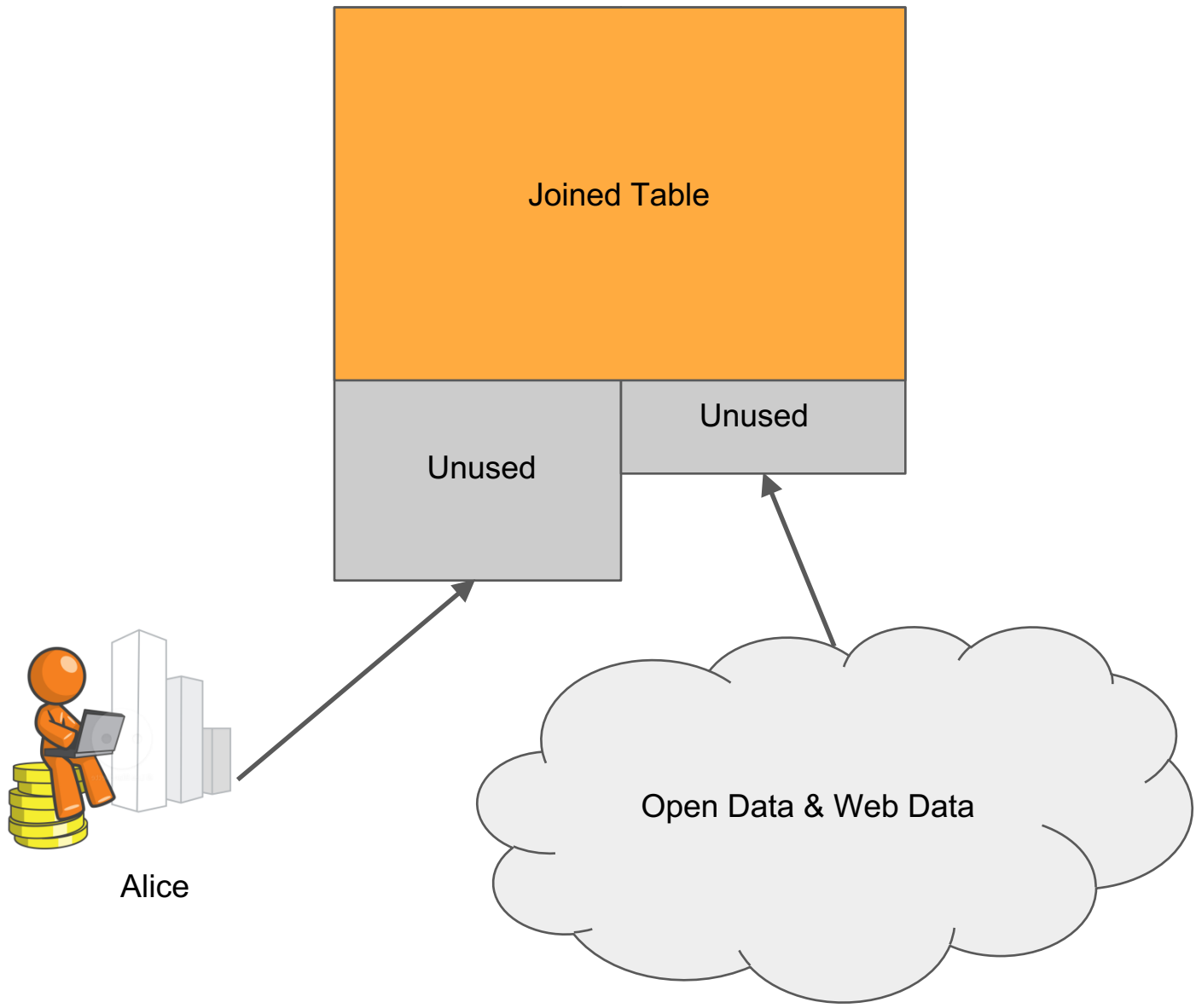


Original Table

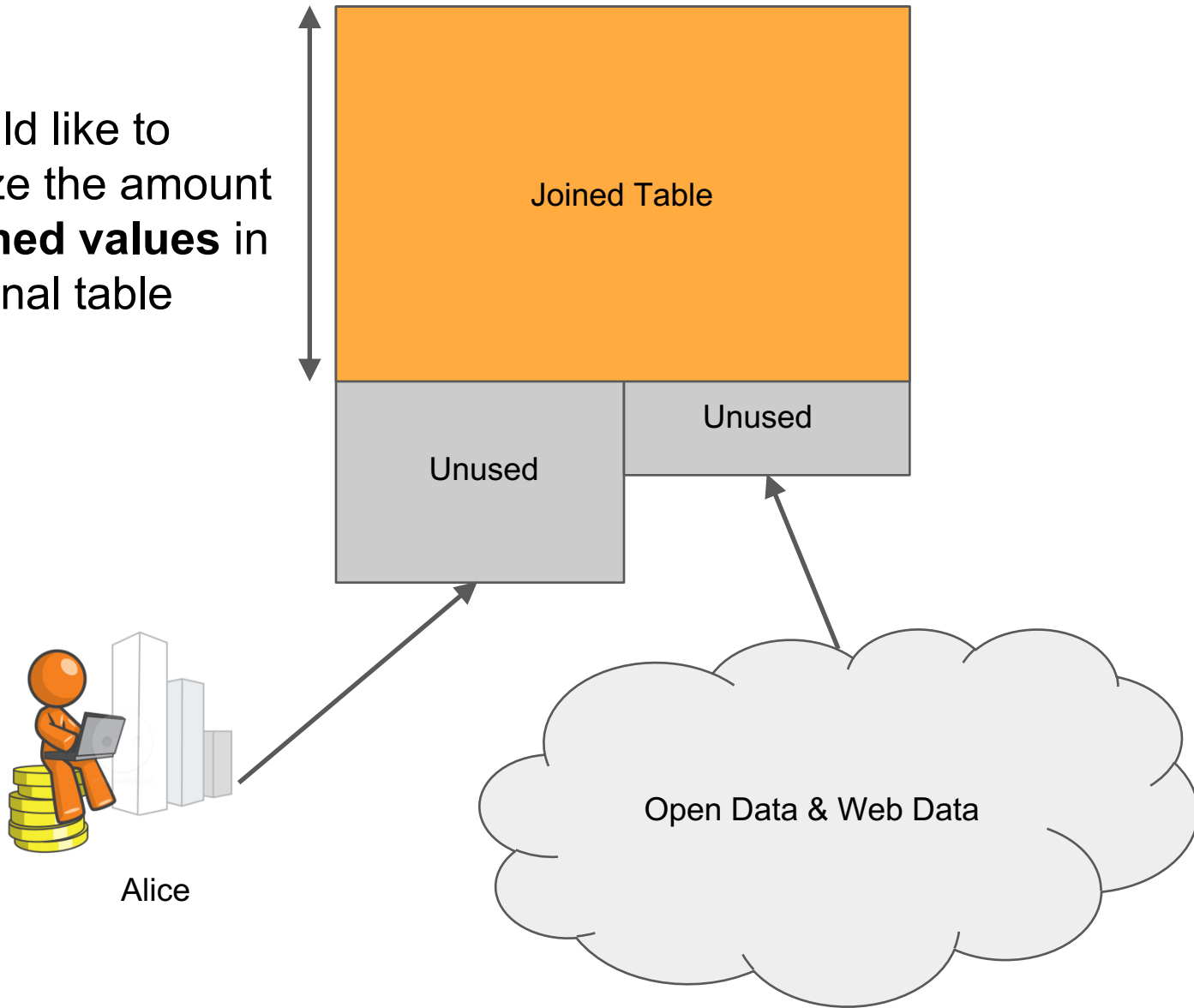
Table Found

Open Data & Web Data

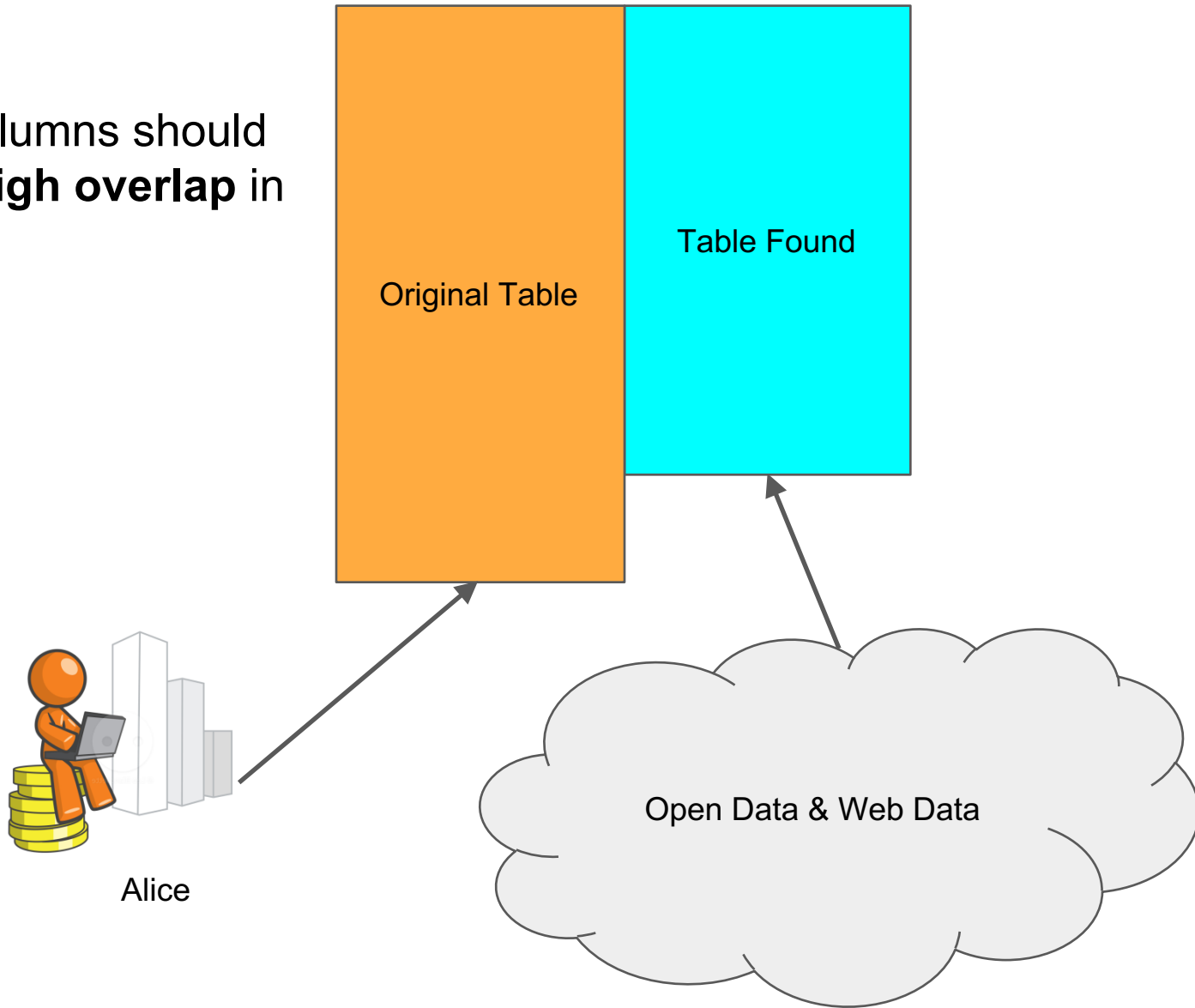
Alice



We would like to maximize the amount of **retained values** in the original table

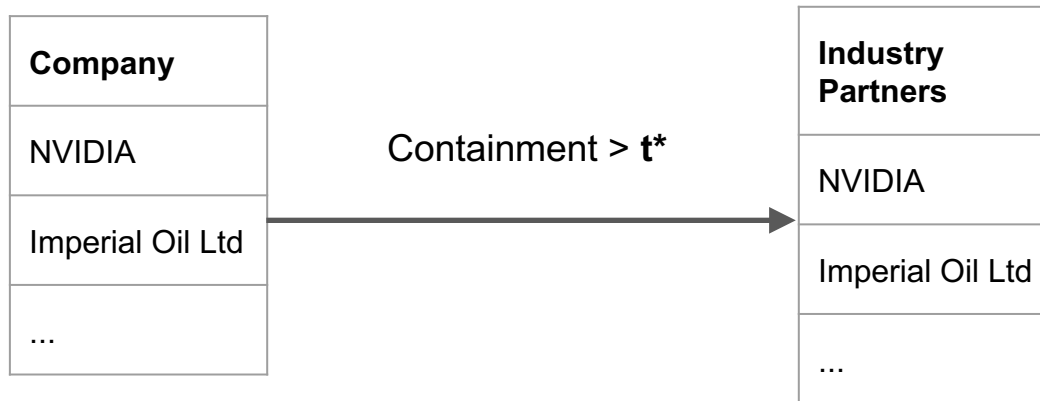


Join columns should have **high overlap** in values



Domain Search

Domain: a set of values in a dataset (e.g. a column in a table)



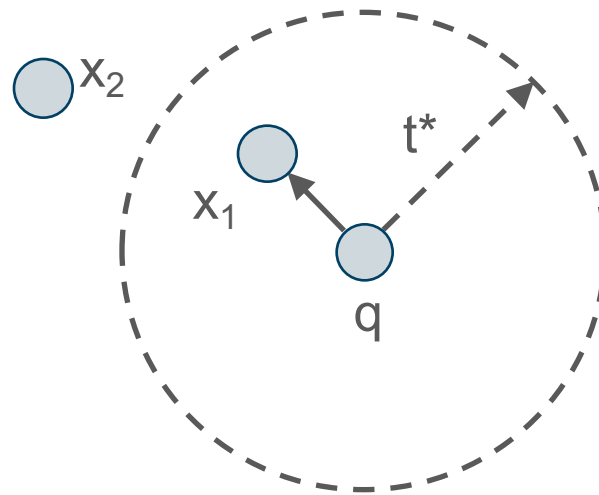
Domain Search: given a query domain Q and **threshold** t^* , find all domains X such that **Containment**(Q, X) is greater than t^* .

$$\text{Containment}(Q, X) = \frac{|Q \cap X|}{|Q|}$$

Approximate Domain Search

We borrow an insight from Approximate Nearest-Neighbour Search.

$$P(x_1) > P(x_2)$$



Existing Solution

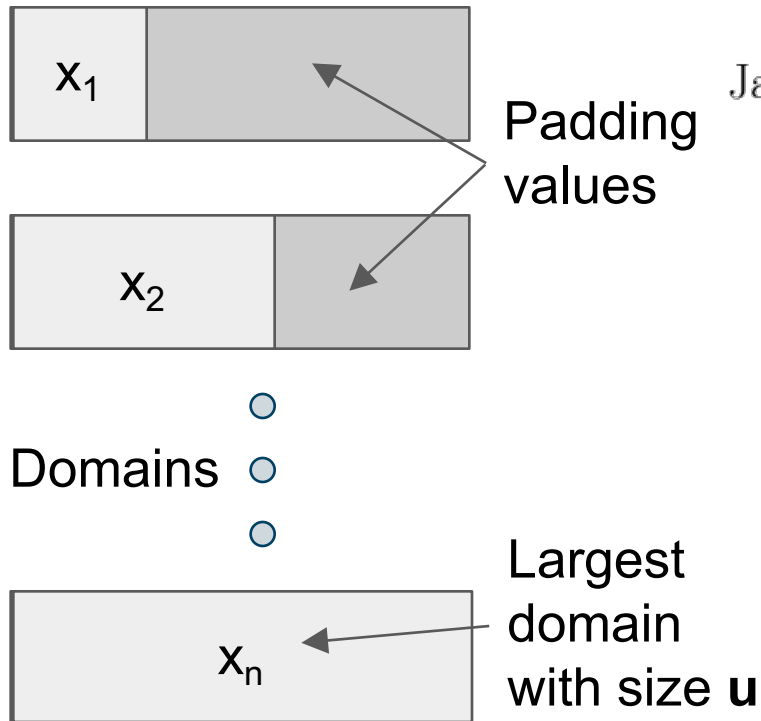
MinHash LSH [Broder 97, Indyk 98] is a search index for Jaccard.

$$\text{Jaccard}(Q, X) = \frac{|Q \cap X|}{|Q \cup X|}$$

- Can be tuned for Jaccard threshold
- Constant space requirement for all domain sizes
- Biased against large domains

Existing Solution

Asymmetric MinHash LSH is a search index for containment
[Shrivastava 15].



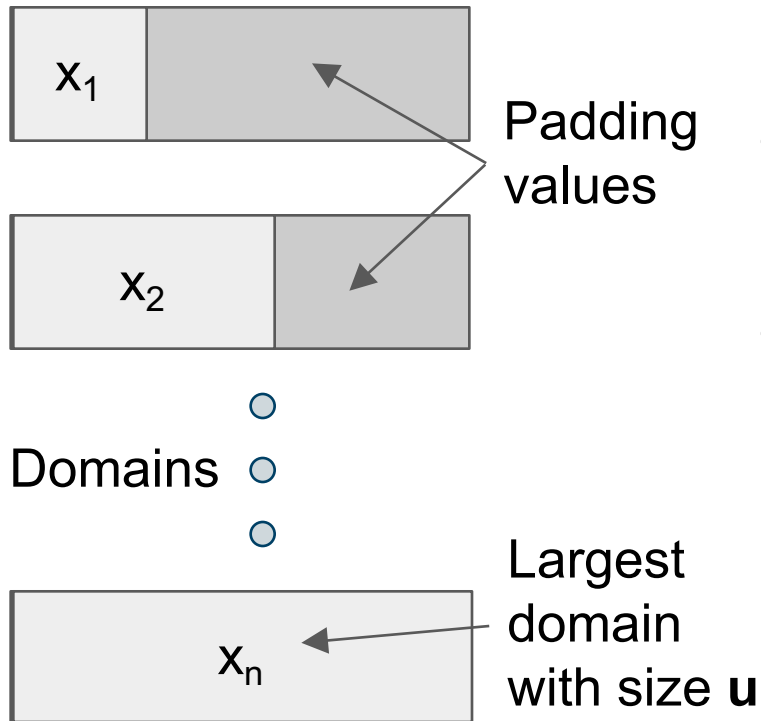
$$\text{Jaccard}(Q, X_{padded}) = \frac{\text{Containment}(Q, X)}{\frac{u}{|Q|} + 1 - \text{Containment}(Q, X)}$$

$$\text{Containment}(Q, X) \propto \text{Jaccard}(Q, X_{padded})$$

MinHash LSH is used to
index the padded domains

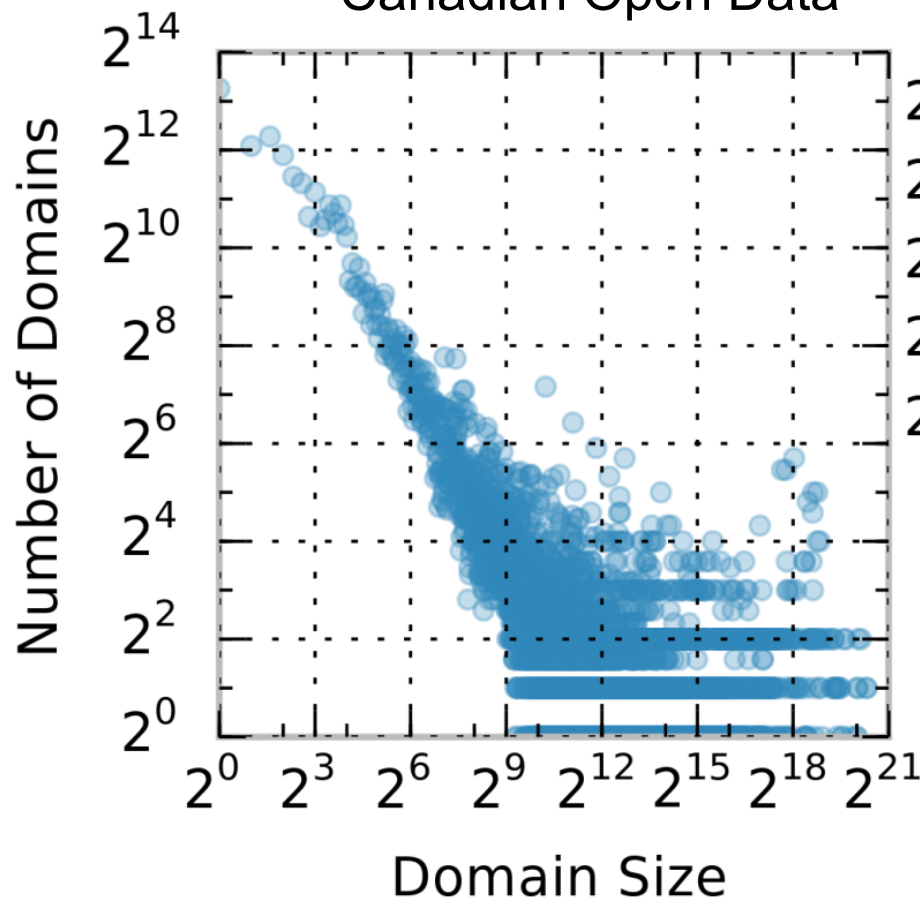
Existing Solution

Asymmetric MinHash LSH is a search index for containment [Shrivastava 15].

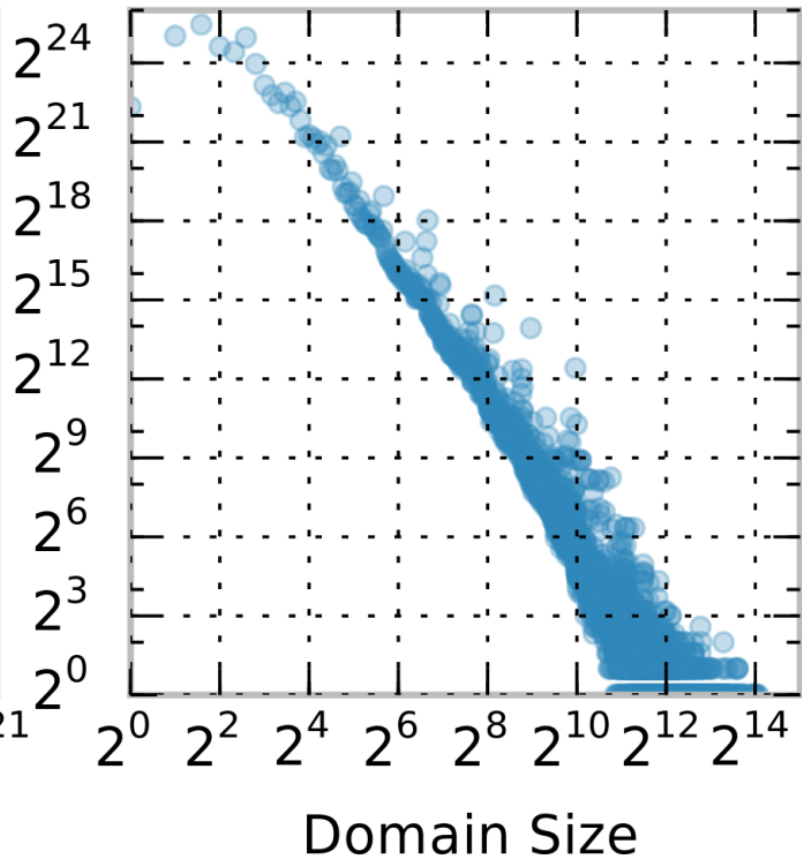


- Difficult to tune for containment threshold
- Padding reduces MinHash accuracy, especially when domain sizes have a skewed distribution

Canadian Open Data



WDC Web Table

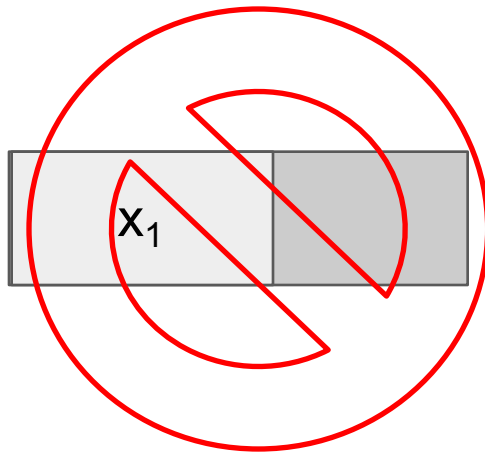


Research Gap: *we need an index for containment search that maintains high accuracy on **skewed** domain size distribution and is tunable for containment **threshold**.*

LSH Ensemble

(Our Contribution)

Can we use MinHash LSH for containment search without padding?



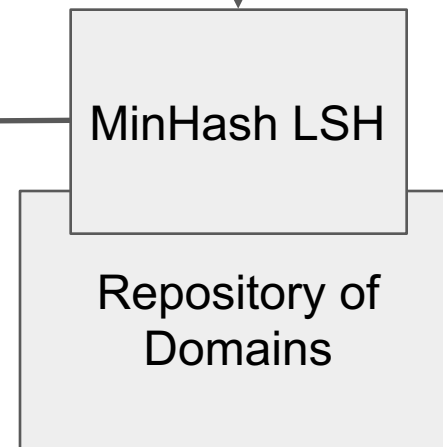
$$\text{Jaccard}(Q, X) = \frac{\text{Containment}(Q, X)}{\frac{|X|}{|Q|} + 1 - \text{Containment}(Q, X)}$$

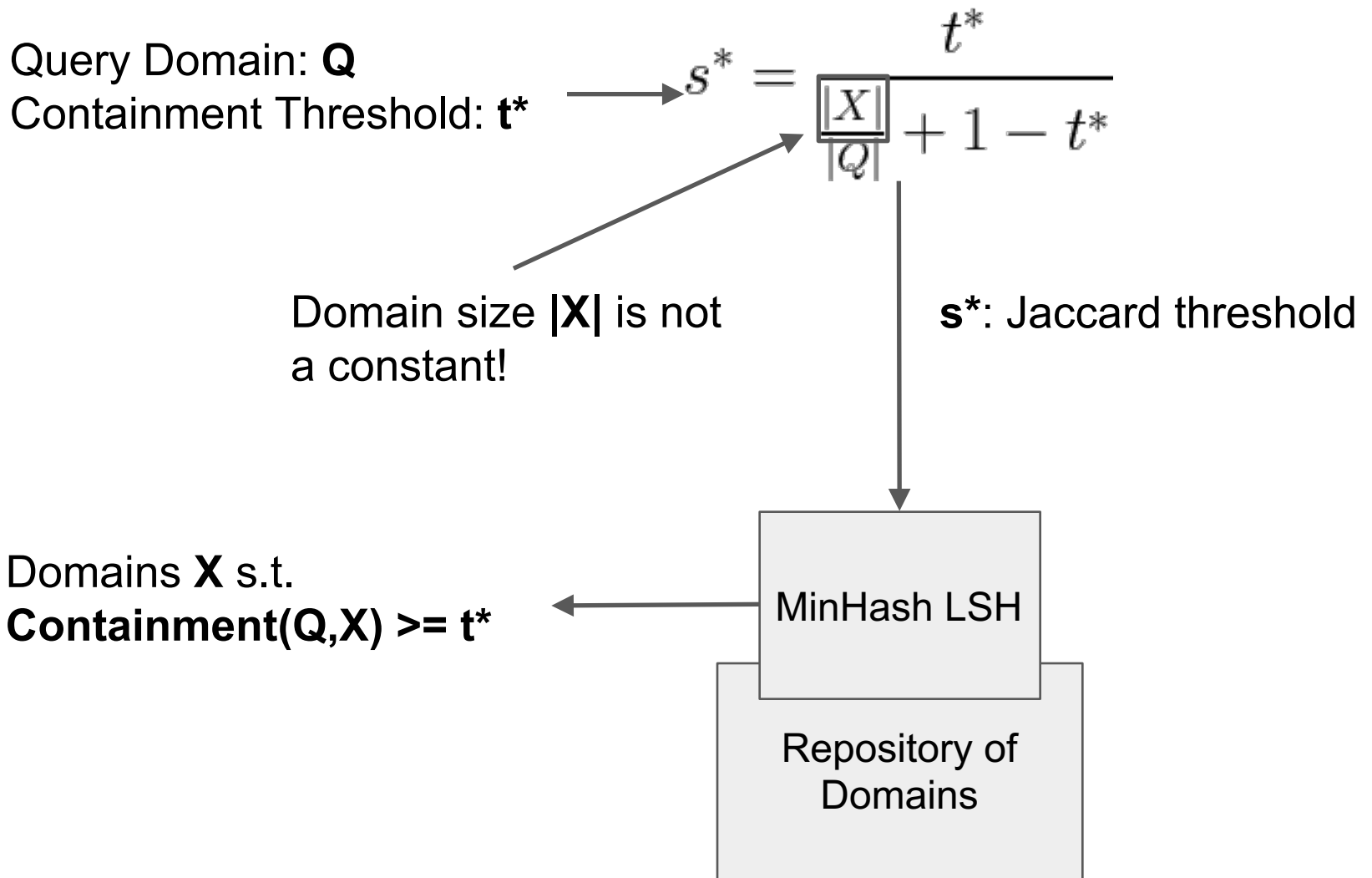
Query Domain: **Q**
Containment Threshold: **t***

$$s^* = \frac{t^*}{\frac{|X|}{|Q|} + 1 - t^*}$$

s*: Jaccard threshold

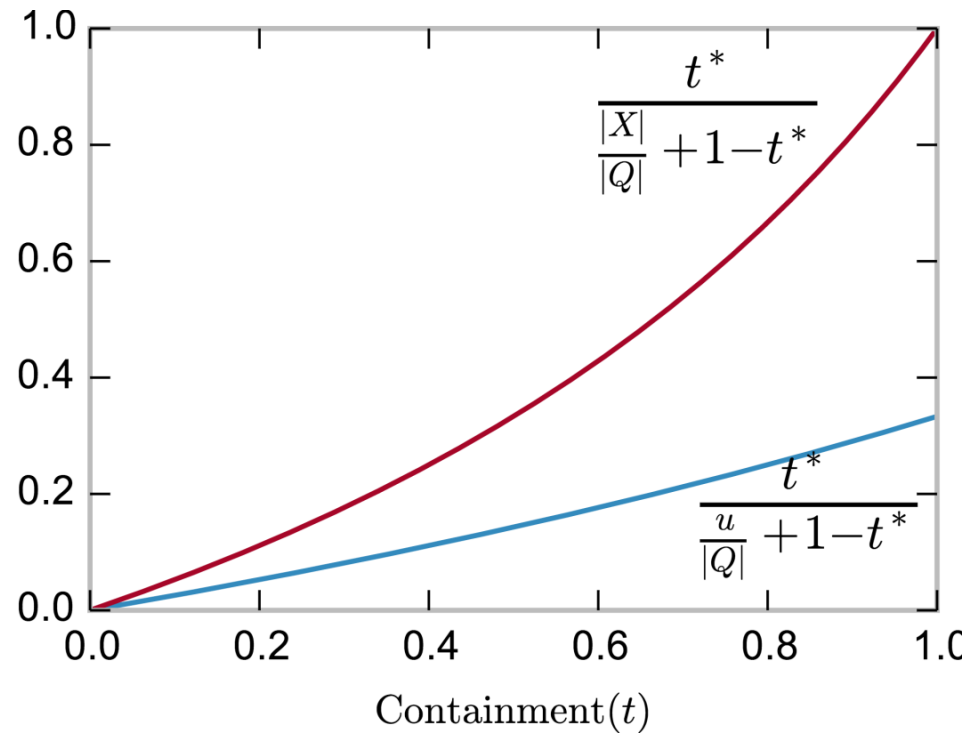
Domains **X** s.t.
Containment(Q,X) >= t*





$$\frac{t^*}{\frac{u}{|Q|} + 1 - t^*} \leq \frac{t^*}{\frac{|X|}{|Q|} + 1 - t^*}$$

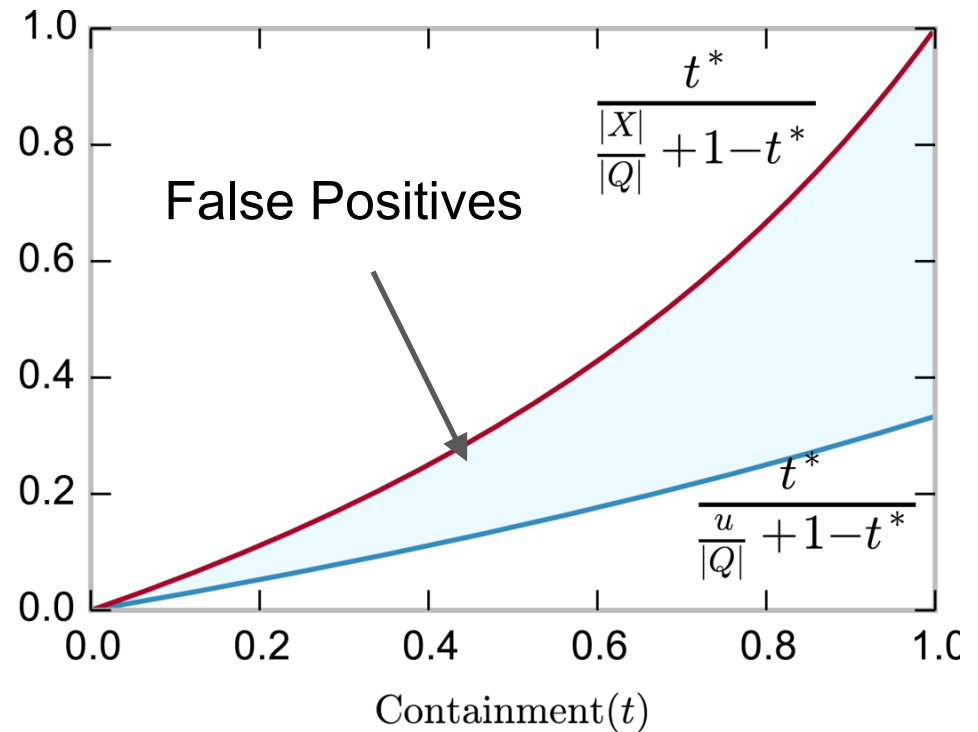
u: the upper bound of domain sizes in range **[l, u]**



$$s_u^* = \frac{t^*}{\frac{u}{|Q|} + 1 - t^*}$$

The new threshold introduces

- false positive domains
- no false negative domains

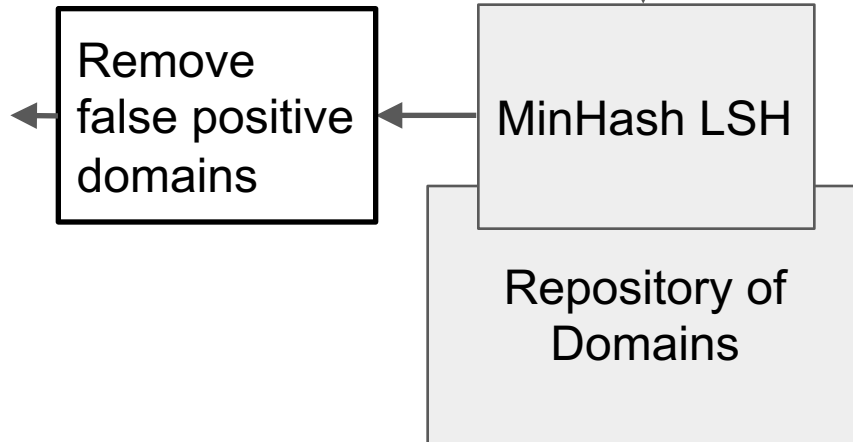


Query Domain: **Q**
Containment Threshold: **t***

$$s_u^* = \frac{t^*}{\frac{u}{|Q|} + 1 - t^*}$$

s_u^* : New Jaccard threshold

Domains **X** s.t.
Containment(Q,X) $\geq t^*$



Query Domain: \mathbf{Q}

Containment Threshold: $\mathbf{t^*}$

$$\longrightarrow s_u^* = \frac{t^*}{\frac{u}{|Q|} + 1 - t^*}$$

s_u^* : New Jaccard threshold

Domains \mathbf{X} s.t.

Containment(Q,X) $\geq t^*$

Remove
false positive
domains

MinHash LSH

Repository of
Domains

Additional cost

The query cost can be reduced if we produce less false positive domains.

$$T_{\text{containment}} = T_{\text{Jaccard}} + \Theta(\text{correct domains}) + \Theta(N_{l,u}^{FP})$$

Number of false positive domains in the range **[l, u]**

Number of
domains in $[l, u]$



$$N_{l,u}^{FP} \leq N_{l,u} \cdot \frac{u - l + 1}{2u}$$

Tight upper bound for
number of false positives
in range $[l, u]$

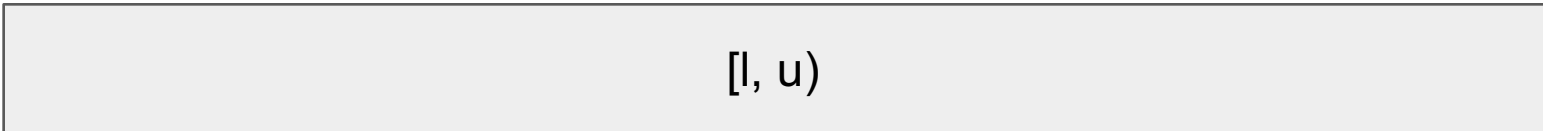
Number of domains in $[l, u]$

We can reduce the range size to reduce the number of false positives domains.

$$N_{l,u}^{FP} \leq N_{l,u} \cdot \frac{u - l + 1}{2u}$$

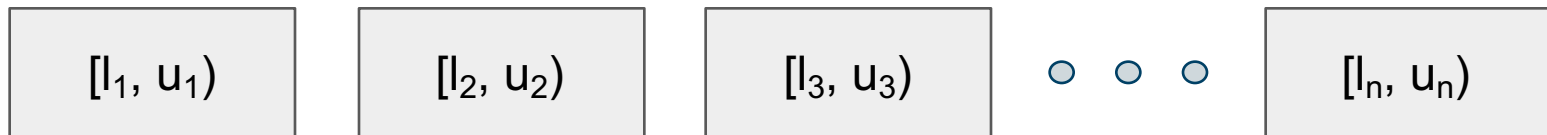
Tight upper bound for number of false positives in range $[l, u]$

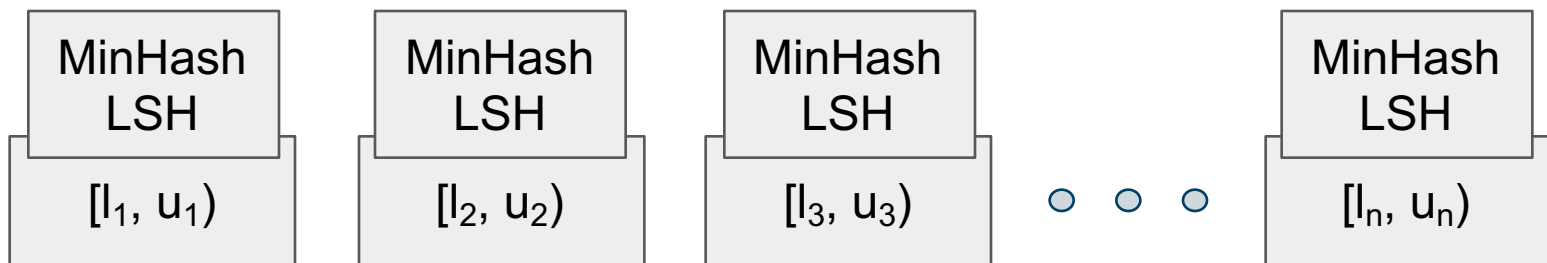
Sorted by domain size



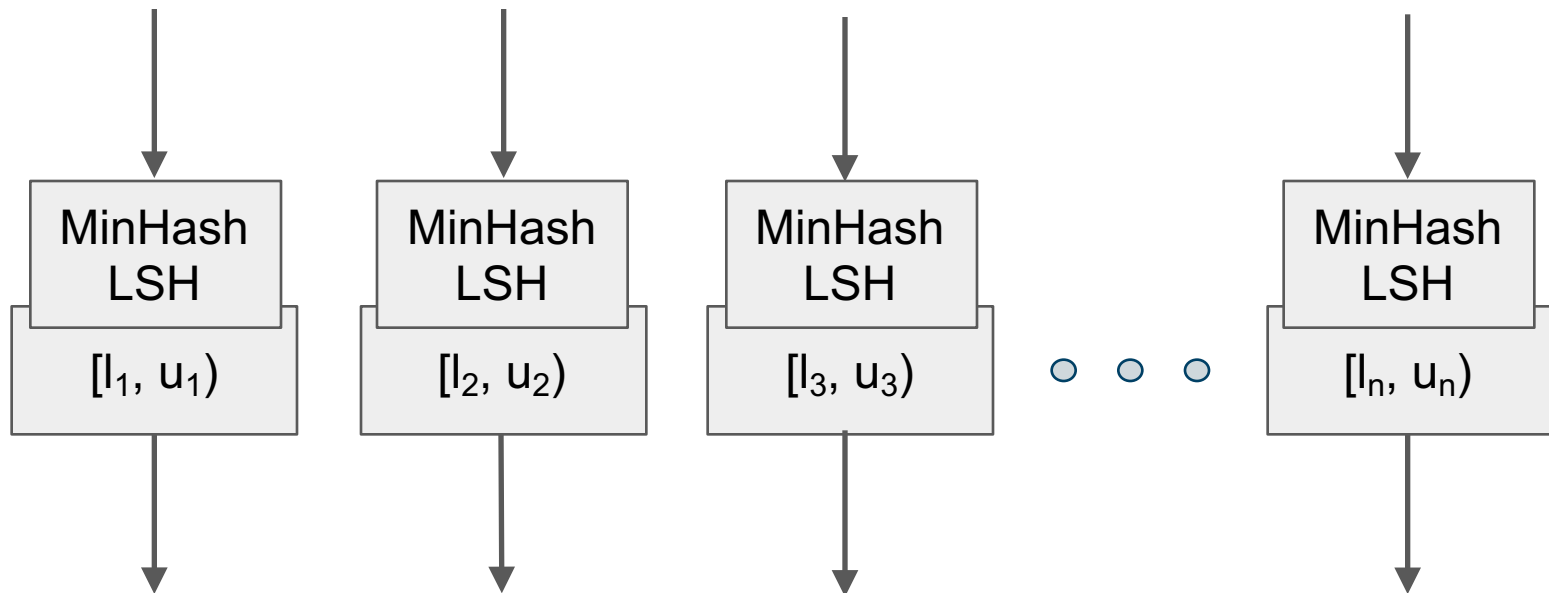
$[l, u)$

Contiguous domain **partitions**





Execute query in parallel



Query cost is determined by the partition with the most false positive domains

This led us to formulate an **optimization problem** for partitioning using the upper bound of $\mathbf{N}_{l,u}^{FP}$ on each partition.

$$\Pi^* = \arg \min_{\Pi} \max_{1 \leq i \leq n} M_i$$

$$M_i = N_{l_i, u_i} \cdot \frac{u - l + 1}{2u}$$

This is equivalent to finding a partitioning such that all partitions have the same \mathbf{M}_i .

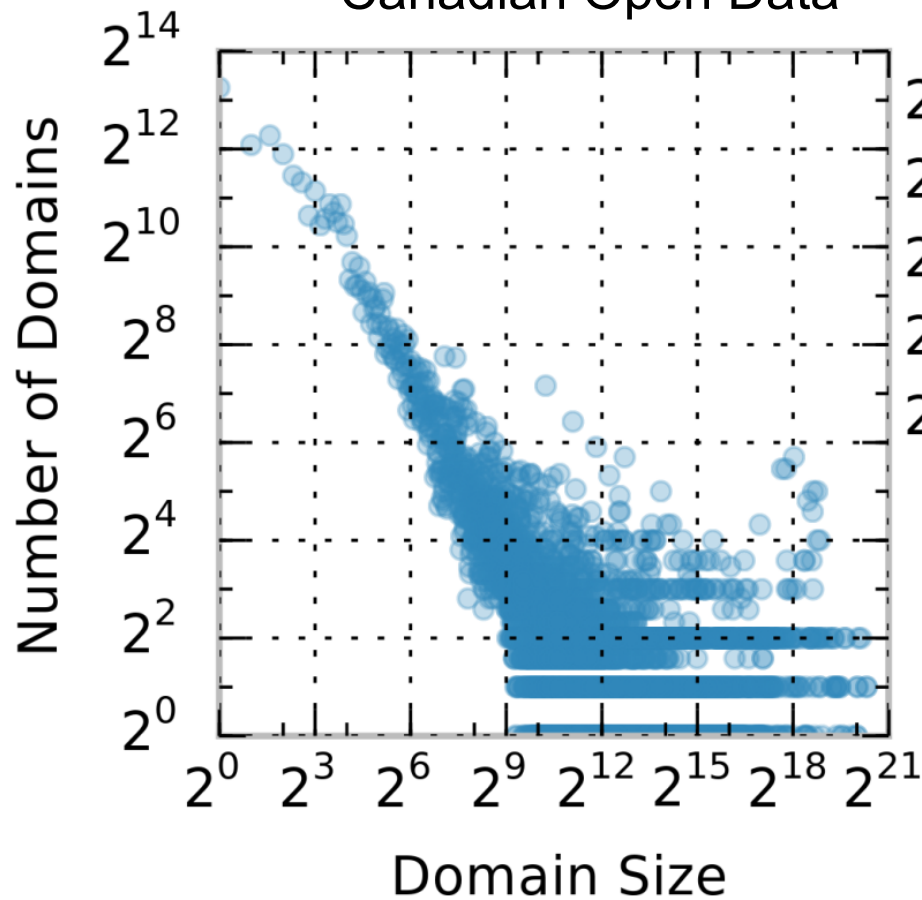
$$\exists \Pi^* \text{ s.t. } M_i = M_j, \forall i, j$$

So far, we have shown:

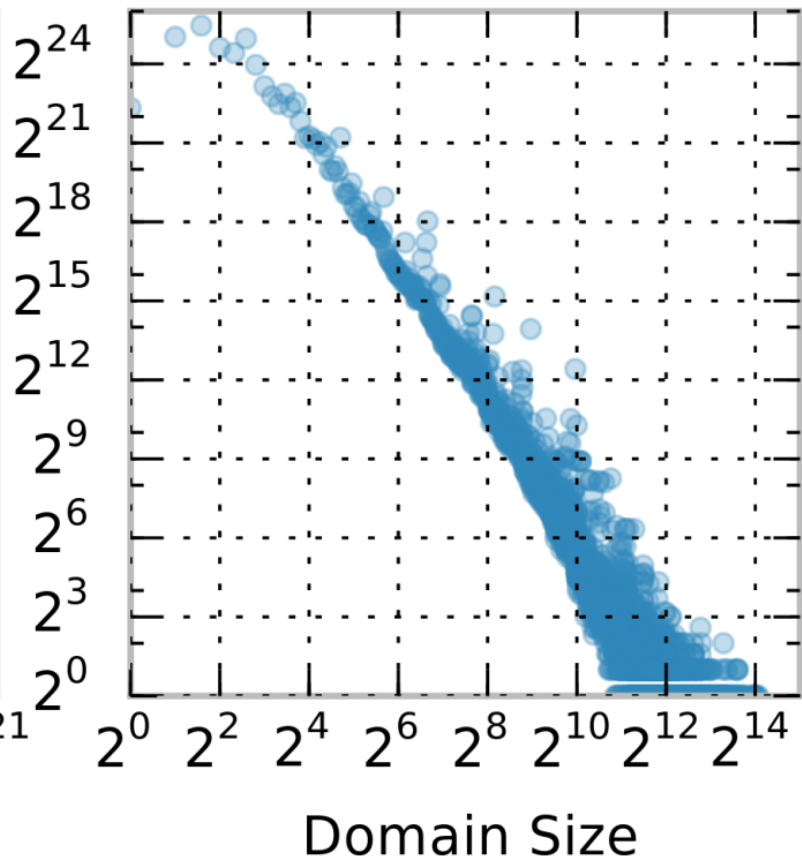
- 1. Partitioning improves query cost by reducing false positives, while maintaining accuracy*
- 2. An optimal partitioning can be verified using a closed form equation*

An optimal partitioning for a real-world domain size distribution?

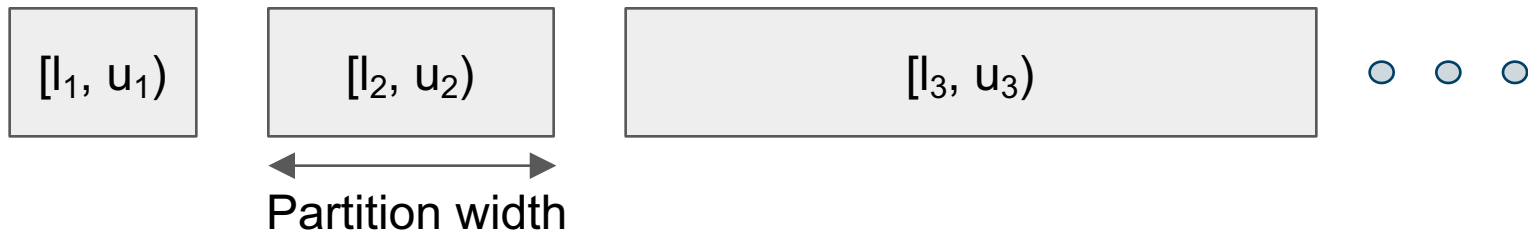
Canadian Open Data



WDC Web Table

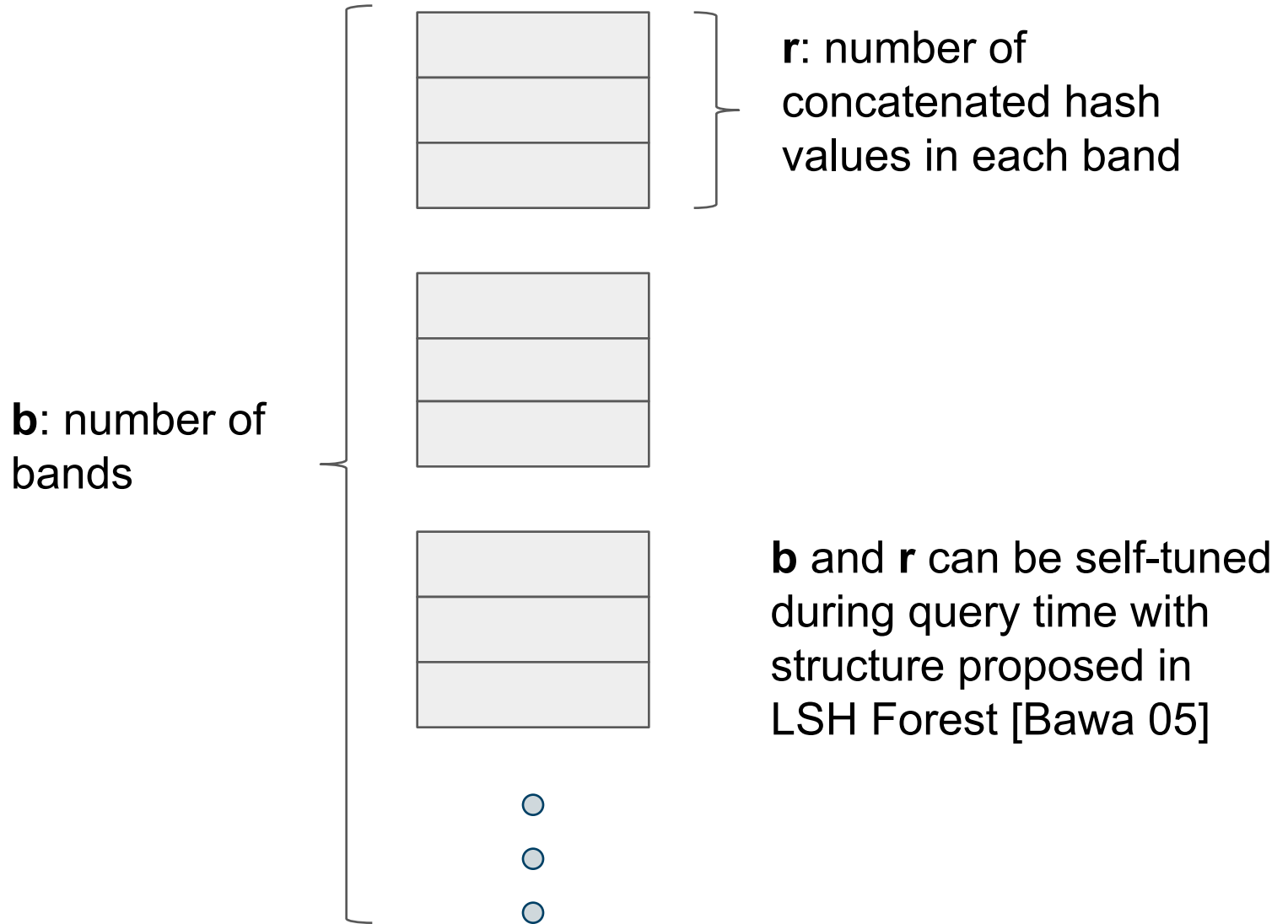


We proved **equi-depth** partitioning is optimal for domains following **power-law** distribution

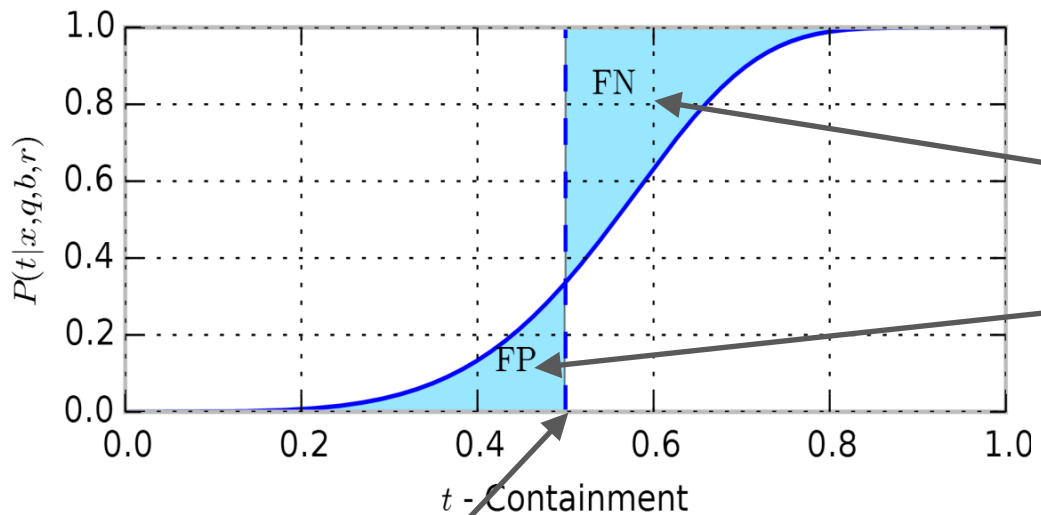


*One last thing: **tune** MinHash LSH for containment threshold*

MinHash LSH [Indyk 98]



We derived the probability of returning with respect to containment given parameters **b** and **r**



Our goal is to minimize the **sum** of false positive and negative probabilities

Containment threshold

Experimental Result

Compared against Asym MinHash [Shrivastava 2015] and MinHash LSH (using our Containment-to-Jaccard conversion) on accuracy, using Canadian Open Data (65,533 domains):

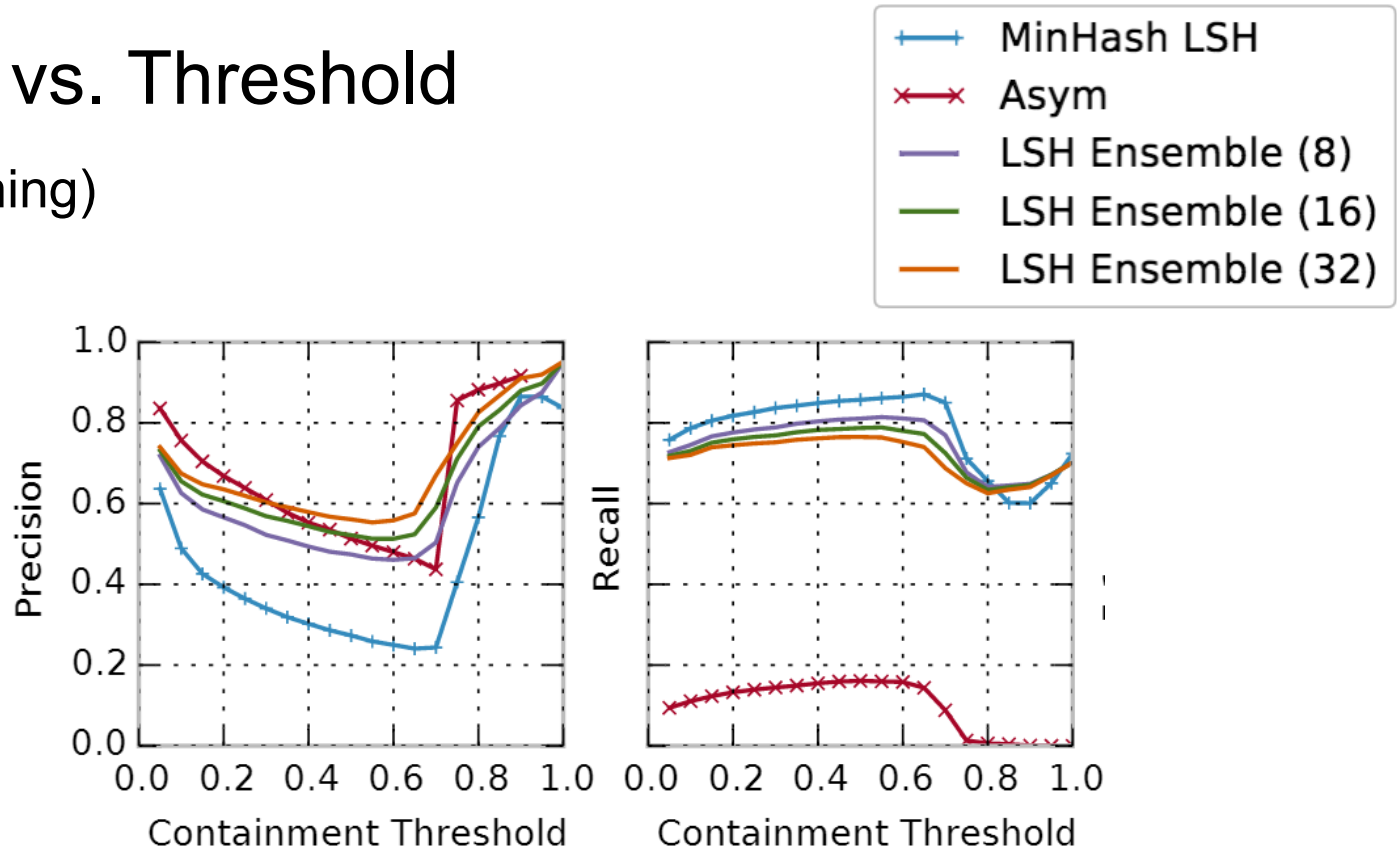
- LSH Ensemble consistently out-performs other techniques
- **More partitions leads to higher accuracy before pruning false positives**

Performance experiment used the complete 2015 WDC English Relational Web Table corpora (263 million domains):

- Mean query time around 3 seconds at 32 partitions
- **More partitions leads to lower query cost**

Accuracy vs. Threshold

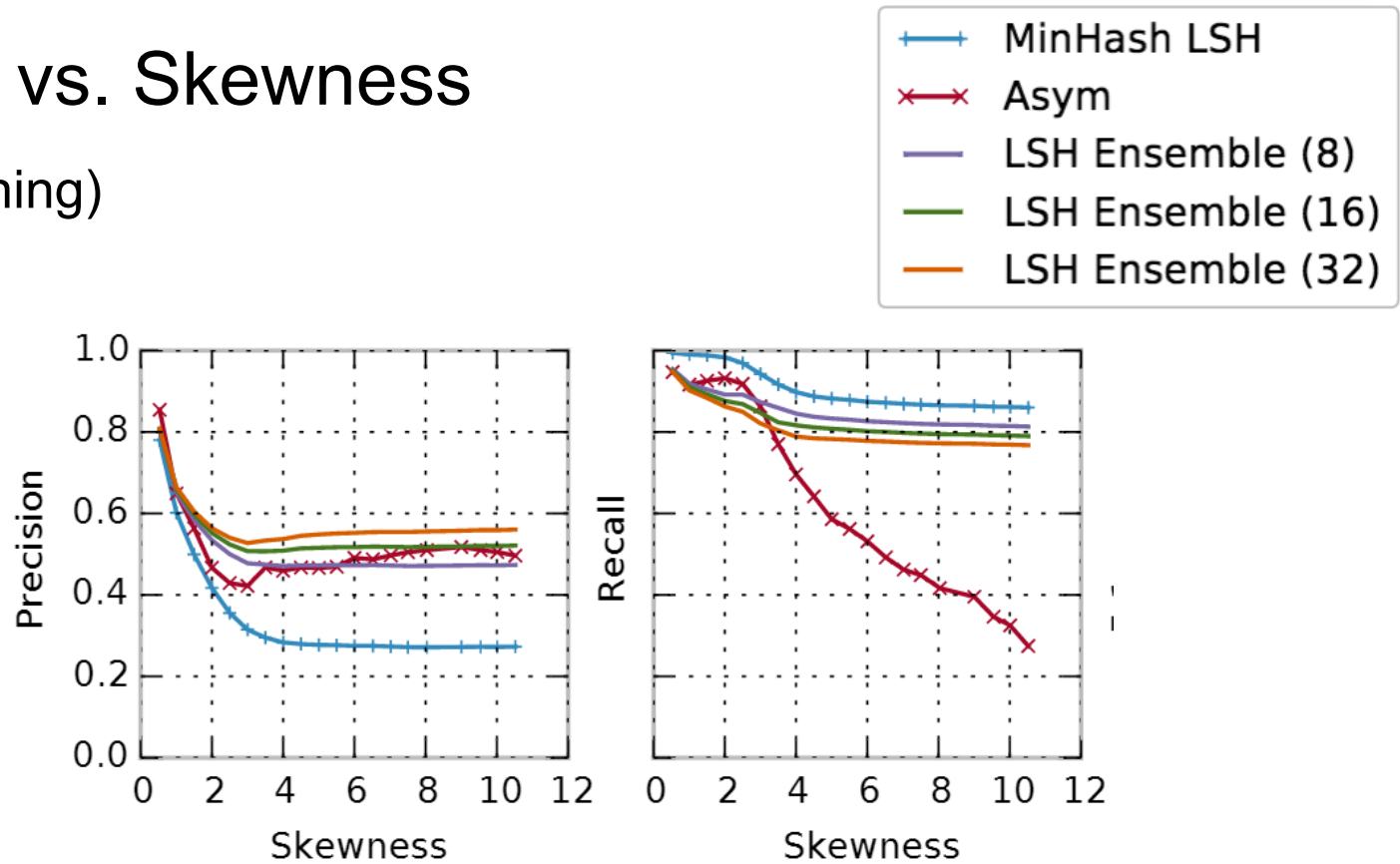
(Before Pruning)



- Creating more partition leads to fewer false positives, while maintaining recall
- Asymmetric MinHash LSH [Shrivastava 15] has high precision, but low recall due to padding

Accuracy vs. Skewness

(Before Pruning)



- Skewness in domain sizes have negative impact on accuracy for all indices
- LSH Ensemble handles skewness better than others

Query Performance

	Mean Query (sec)	Precision Before Pruning ($t^*=0.5$)
MinHash LSH	45.13	0.27
LSH Ensemble (8)	7.55	0.48
LSH Ensemble (16)	4.26	0.53
LSH Ensemble (32)	3.12	0.58

On 263 million domains
(WDC Web Table)

Speed up is due to:

- Fewer false positive domains to process (higher precision)
- Parallelization

Recap

LSH Ensemble

- Uses multiple MinHash LSH built on domain size partitions to approximate containment search and maintain accuracy
- Optimal partitioning (equi-depth) for power-law distributions
- Self-tunable at query time given any threshold

Thank you!

Erkang (Eric) Zhu
ekzhu@cs.toronto.edu