

Using AI, Mathematics, and Statistics to Find Similar Data in Massive Data Ecosystems

Collibra

Eric Warner, PhD
Senior Manager, AI Engineering

What is Collibra? | Collibra is the system of engagement for data



Collibra Data Intelligence Platform

Enterprise

Data Catalog

Governance

Lineage

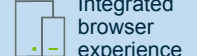
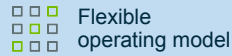
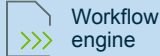
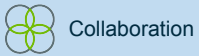
Quality & Observability

Data Marketplace

Privacy and Security

Active metadata graph

Core services



Data Stores and Applications



Hyperscalers



AI and BI Tools



Integrations

100+
pre-built integrations

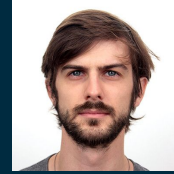
Who? | Collibra Artificial Intelligence Team



Eric Warner
Senior AI Manager



Larry Hau
Senior Director, Product



Vicky Froyen
Staff Product
Manager

Discovery



Gretel De Paepe
Principal Data Scientist

Delivery



Anna Filipiak
Staff AI Engineer



Kelsey Schuster
Senior AI Engineer



Lillian Neff
AI Engineer



Nick Evers
Senior AI Engineer

Operations



Cristian Caraballo
Senior Ops Eng

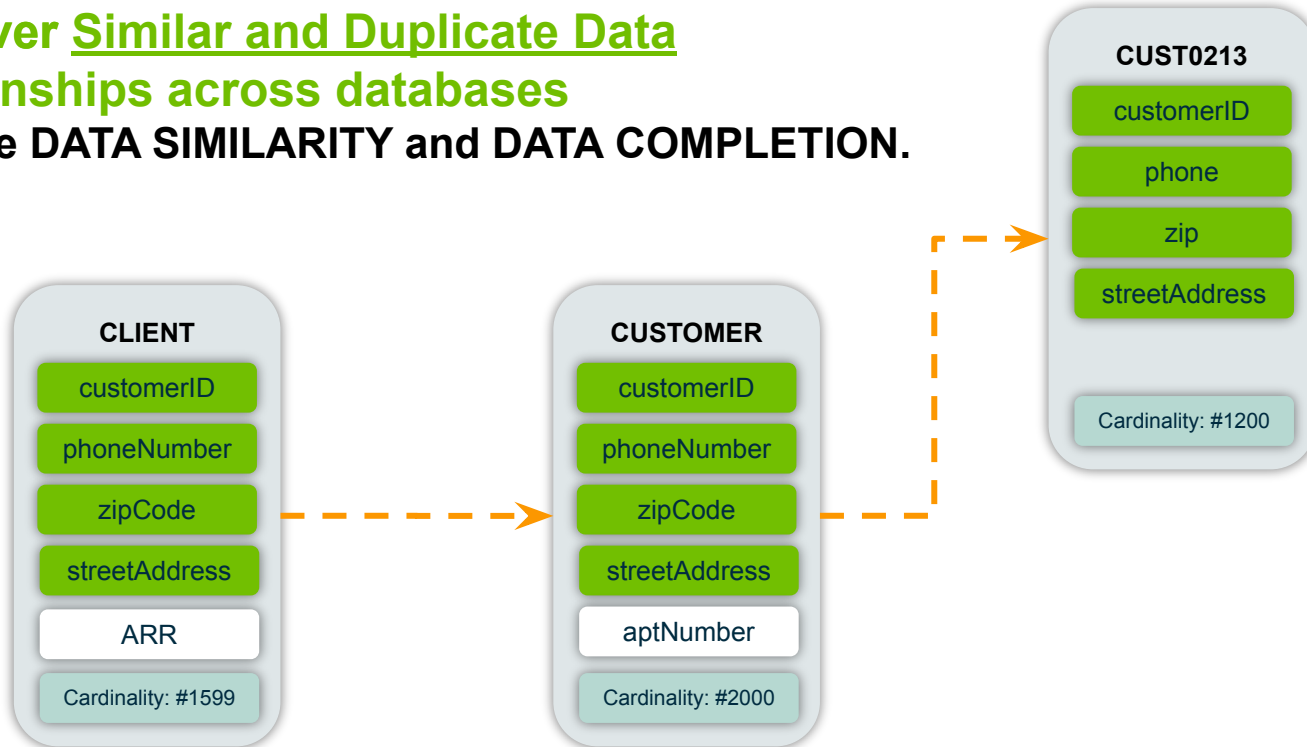


Charlie Ang
Staff Ops Eng

What? | Data Similarity

Discover Similar and Duplicate Data relationships across databases

To drive DATA SIMILARITY and DATA COMPLETION.



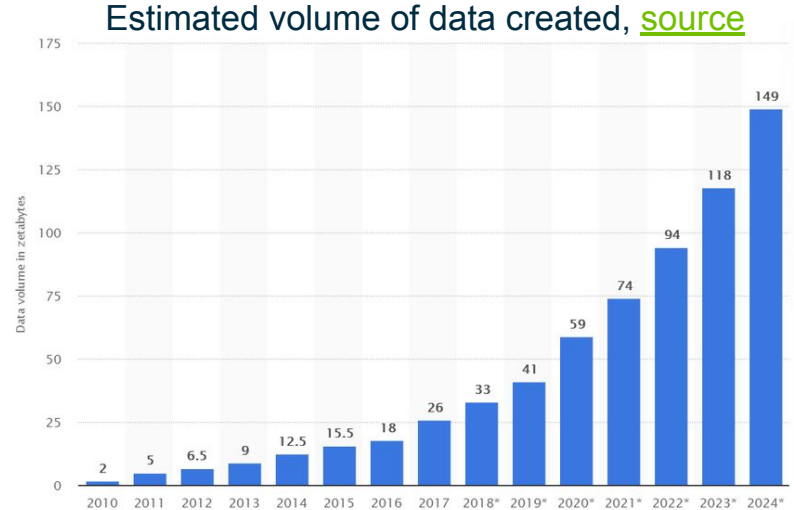
Why? | Data Similarity

Data is a \$273T market¹

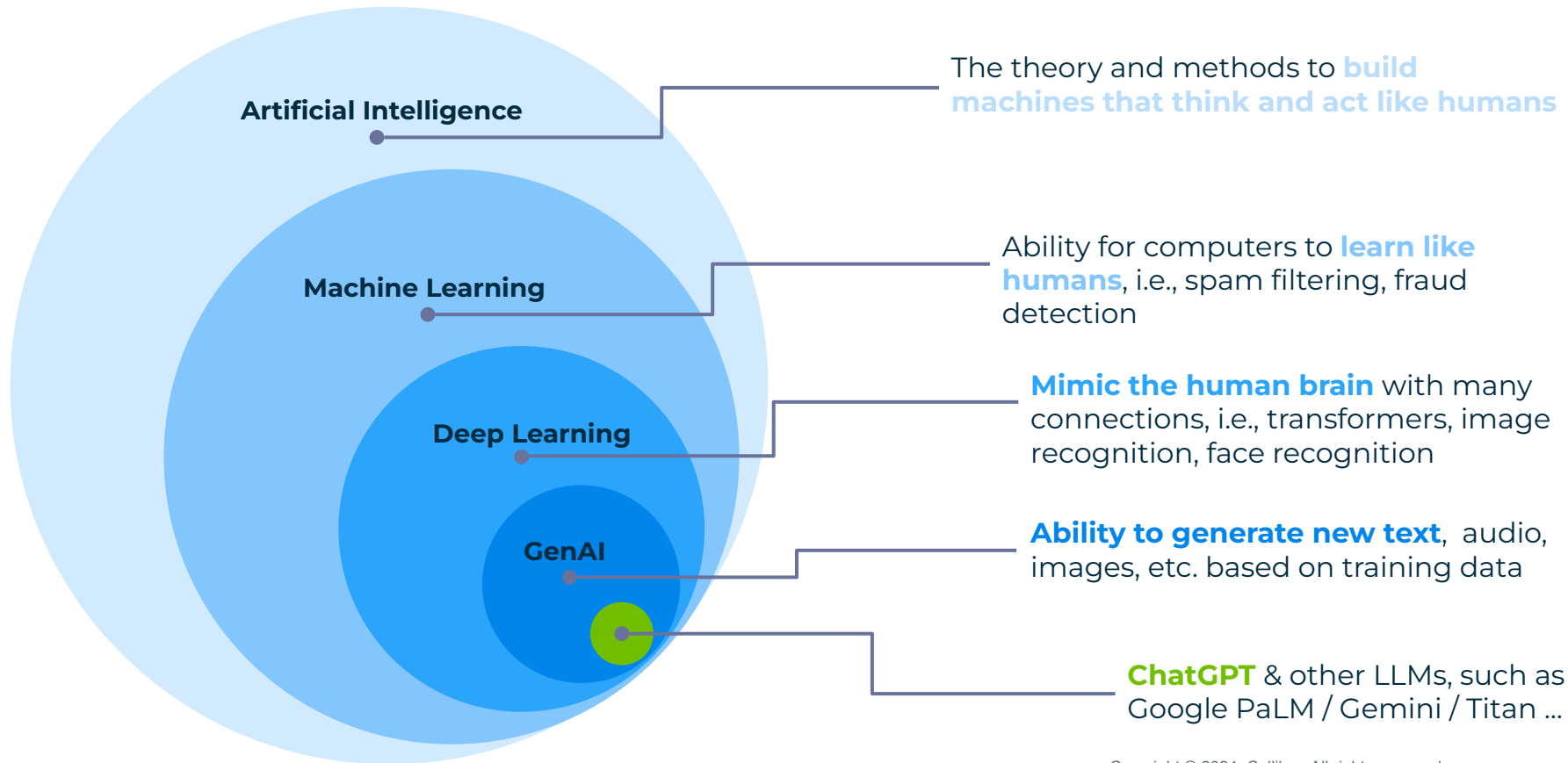
A company with a thousand data workers **wastes time equivalent to \$5.7 million a year trying to find relevant data²**

By 2026, 60% of GenAI applications will fail due to bad, wrong, or not enough data³

^{1,2,3}Source: IDC



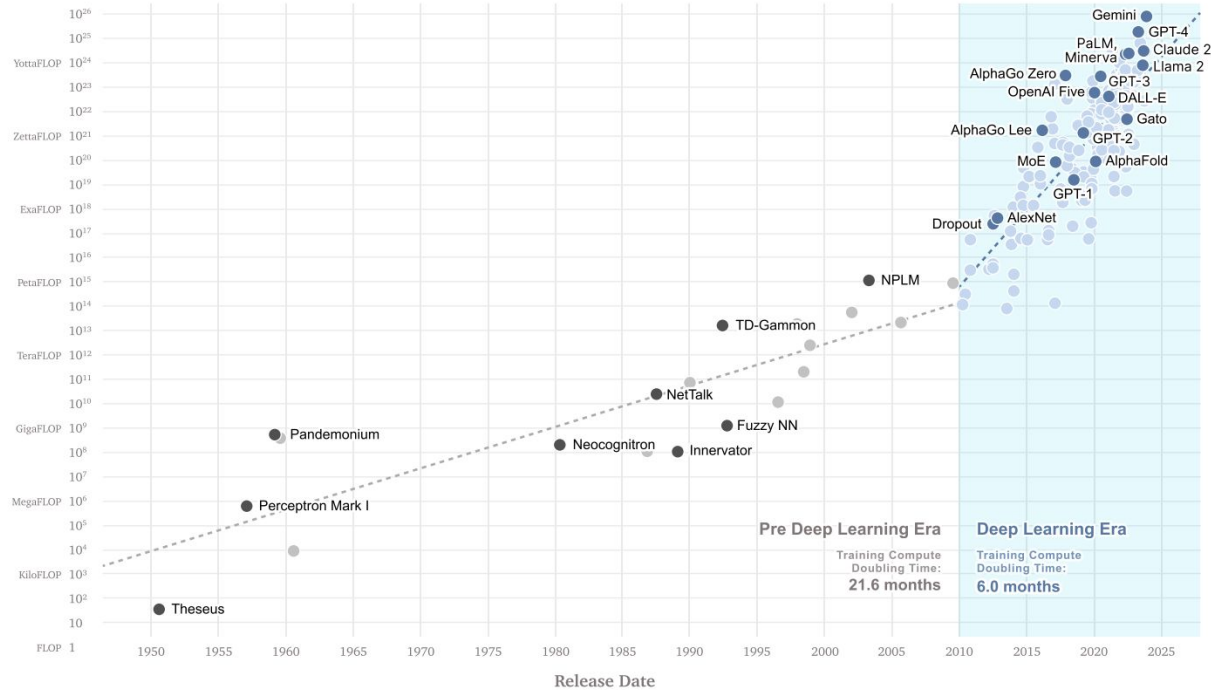
AI | Riding the Tidal Wave of Hype



LLMs do not come without cost...

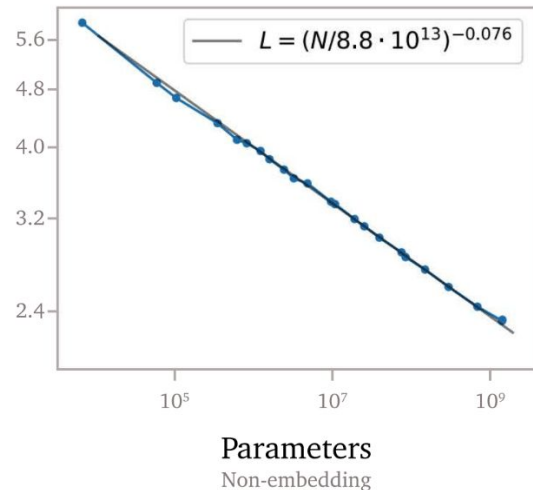
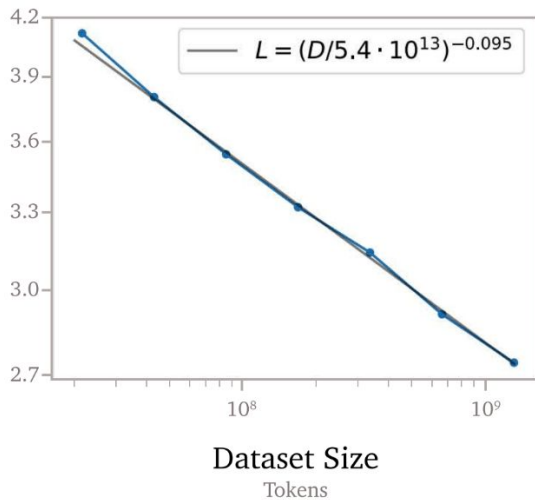
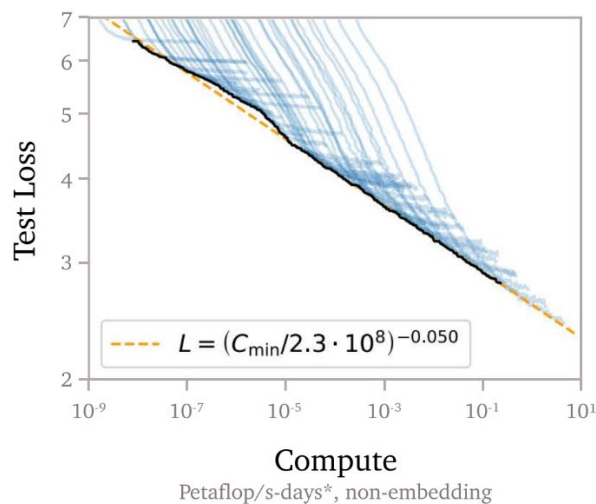
Compute Used for AI Training Runs

Total compute used to train notable AI models, measured in total FLOP (floating-point operations) | Logarithmic



[Source](#)

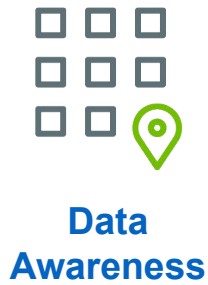
...or upside



* A Petaflop/s-days is equivalent to approximately 10^{20} floating point operations

[Source](#)

Navigating the Keys to Data Science



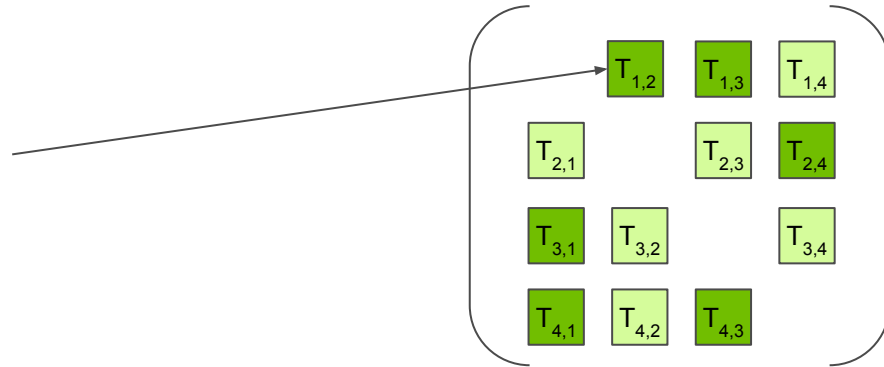
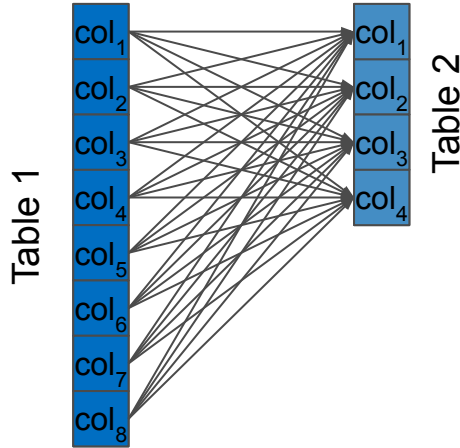
At a high level, how could this (inefficiently) work?



Algorithm
Selection

Comparing an average of n columns per table and
 p entries per column...

... for m Tables



Compute Efficiency: $O(p^2 m^2 n^4 m^2)^*$

Storage Requirements: 300.02 TB

(100,000 tables, 100 cols per table, 50 entries per column)

**250 quadrillion computations in the above example or ~79 years if
GPU can do 100 million comparisons per second*

Clearly, this is not feasible

So how are we actually achieving this?

Let's start at the outcome

Constraints Drive Solutions



Directed
Use-Case

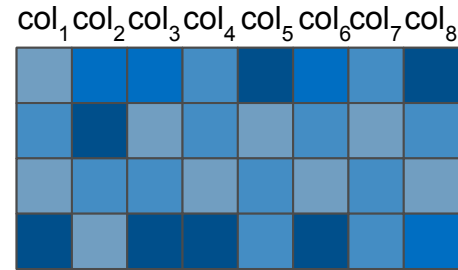
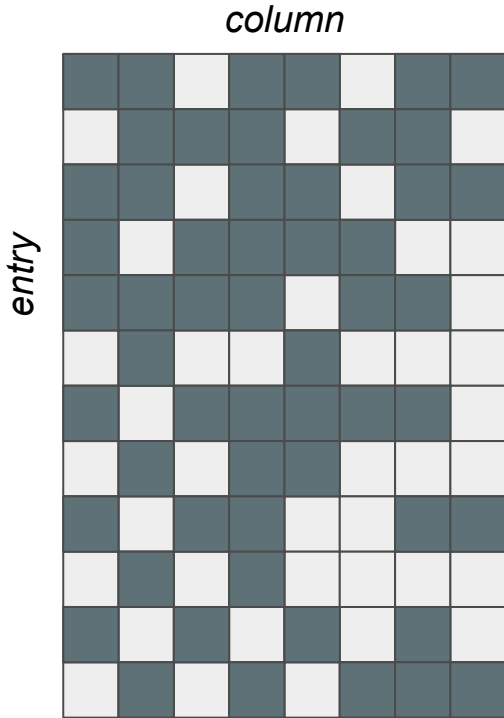
Challenge | Our algorithm

- Must run on customer hardware
- Cannot phone home to an internal or 3rd-party API
- Must actually be able to determine similar/duplicate data

Data Source Constraints



Data
Awareness



...we must compress...

...and add noise...

..while achieving >90% accuracy. How??

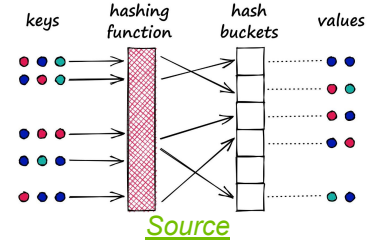
Given a series of columns...
...randomly sampled and streamed in
batches...

Compression | Our Secret Sauce to Performance

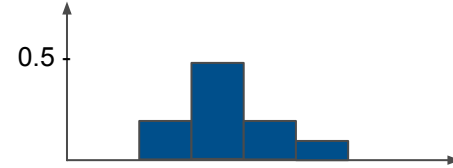


Algorithm Selection

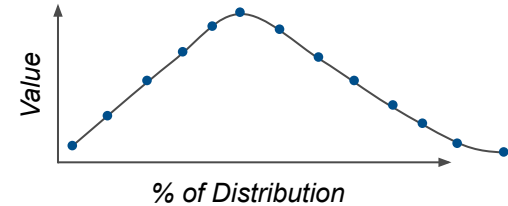
String/Date/Categorical Data
 {'MA', 'CT', 'MI', 'CO', 'MI', 'NH'}



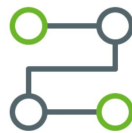
Integer Data
 {1, 2, 2, 2, 3, 2, 2, 3, 1, 4}



Float Data
 {-0.63, -0.41, 0.77, 0.08, 0.28, -0.96, 0.23, 0.06, -0.18, -0.07 }



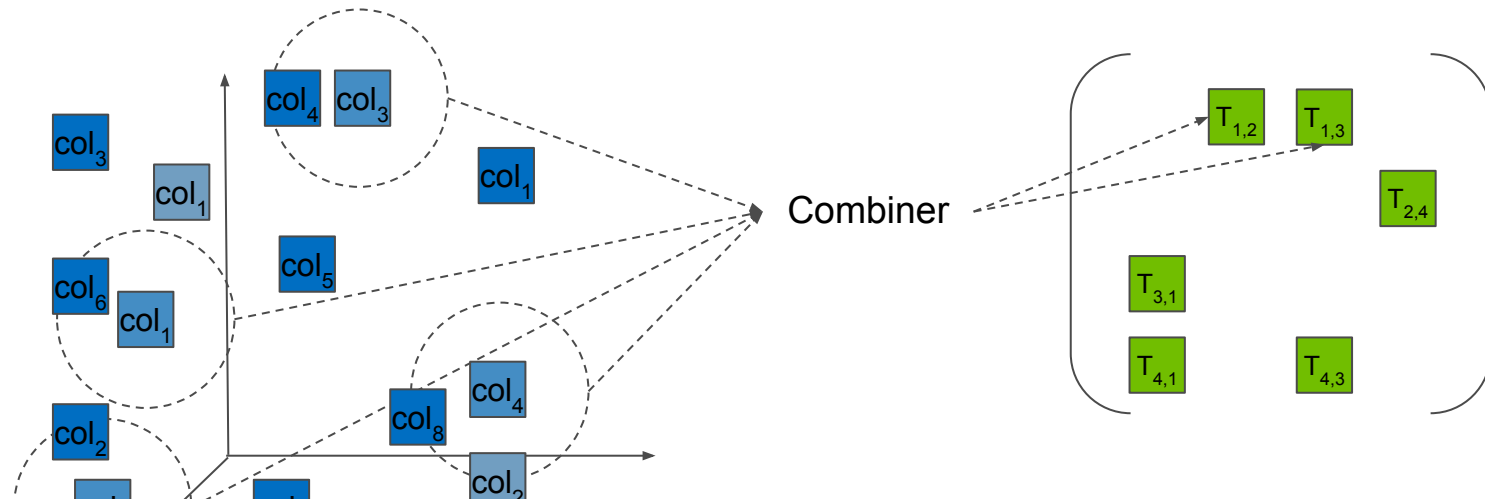
Our schema aims to maintain as much information as possible while maintaining security and legal boundaries



At a high level, how does this actually work?

Comparing an average of n columns per table and p entries per column...

... for m Tables



Compute Efficiency: $O(p * m * n) + m$
 Storage Requirements: **90.002 GB**
 (100,000 tables, 100 cols per table, 50 entries per column) *0.03% of brute-force*

**~100 million computations in the above example if 20% of tables have a duplicate (0.00000004% of brute-force)*

Compression

This is not new....

Morse Code (1840's)

A	· -	N	- ·	1	· - - -	?	· · - -
B	· · · -	O	- - -	2	· · - - -	!	· · - - -
C	· - · -	P	- · - -	3	· - - -	.	· - - -
D	· · · -	Q	- · - ·	4	· - · -	,	- - - -
E	·	R	- · · -	5	· · · ·	;	· - - -
F	· · - ·	S	· · ·	6	· · · ·	:	· - - -
G	- · - ·	T	-	7	- - · -	+	· - - -
H	· · · ·	U	- · -	8	- - - ·	-	· - - -
I	· ·	V	· · -	9	- - - ·	/	· - - -
J	- · - -	W	- · - ·	0	- - - -	=	· - - -
K	- · · -	X	- · - ·				
L	· - · -	Y	- · - ·				
M	- -	Z	- · - ·				

[source](#)

Shannon Coding (1940's)

i	p_i	l_i	$\sum_{n=0}^{i-1} p_n$	Previous value in binary	Codeword for a_i
1	0.36	2	0.0	0.0000	00
2	0.18	3	0.36	0.0101...	010
3	0.18	3	0.54	0.1000...	100
4	0.12	4	0.72	0.1011...	1011
5	0.09	4	0.84	0.1101...	1101
6	0.07	4	0.93	0.1110...	1110

[source](#)

...and very relevant to current trends in AI....



Language Modeling Is Compression

Grégoire Delétang^{*1}, Anian Ruoss^{*1}, Paul-Ambroise Duquenne², Elliot Catt¹, Tim Genewein¹, Christopher Mattern¹, Jordi Grau-Moya¹, Li Kevin Wenliang¹, Matthew Aitchison¹, Laurent Orseau¹, Marcus Hutter¹ and Joel Veness¹

^{*}Equal contributions, ¹Google DeepMind, ²Meta AI & Inria

[source](#)

Burning Question | Similar vs Duplicate?



Directed
Use-Case

This is not a record-to-record match, but rather a **column-to-column statistical comparison**

PROS

- Computational time
- Spatial requirements
- Does not store raw customer data at rest (lower risk)

CONS

- **Sample Accuracy** | Exact row comparisons are not done
- Column comparisons are statistical comparisons, not exact matching
- We see a maximum number of rows

Bottom Line | We believe that this method improves at volume, i.e., with more tables & columns comes more accuracy & utility

Constraints Drive Solutions



Directed
Use-Case

Challenge | Our algorithm

- Must run on customer hardware
- Cannot phone home to an internal or 3rd-party API
- Must actually be able to determine similar/duplicate data

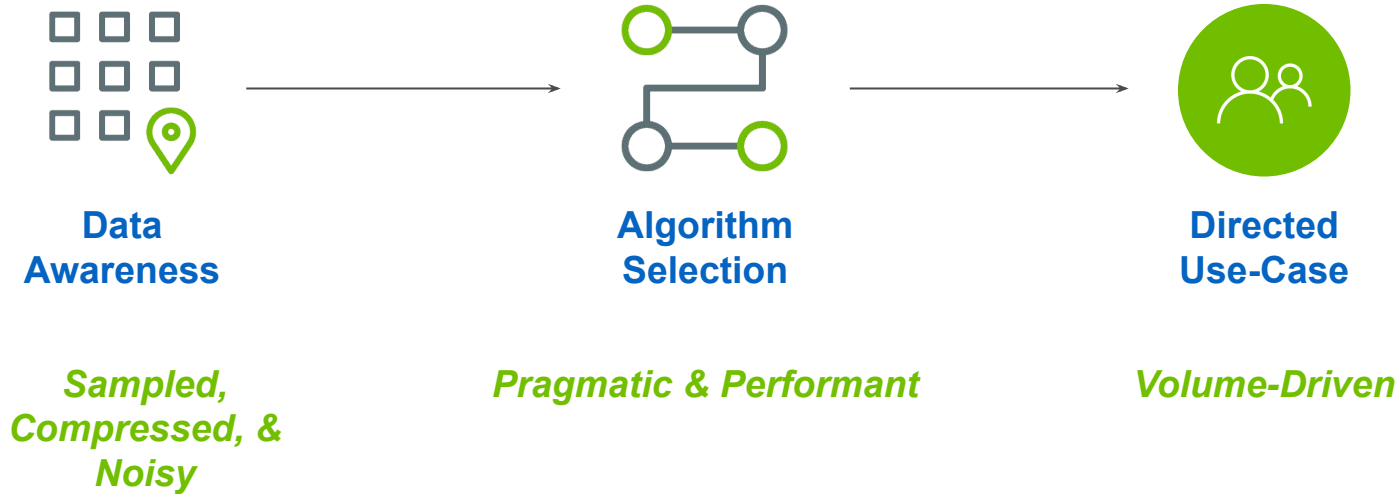
Solution | Let's go back to basics

- Decompose data into a series of statistical properties which minimize compression loss
- Build an algorithm which fits the needs of the end-user



[source](#)

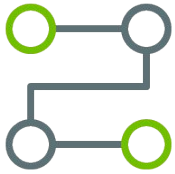
Navigating the Keys to Data Science



Prolonged Value | Technology Solution



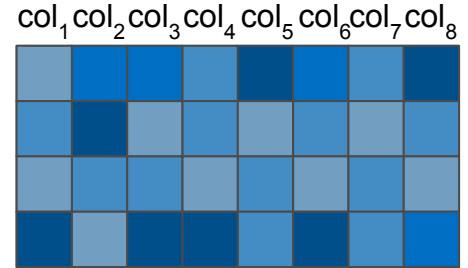
Pipeline to acquire data



Compressed Data Representation



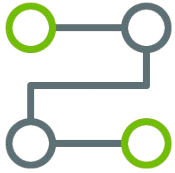
Context surrounding the data



Prolonged Value | Technology Solution



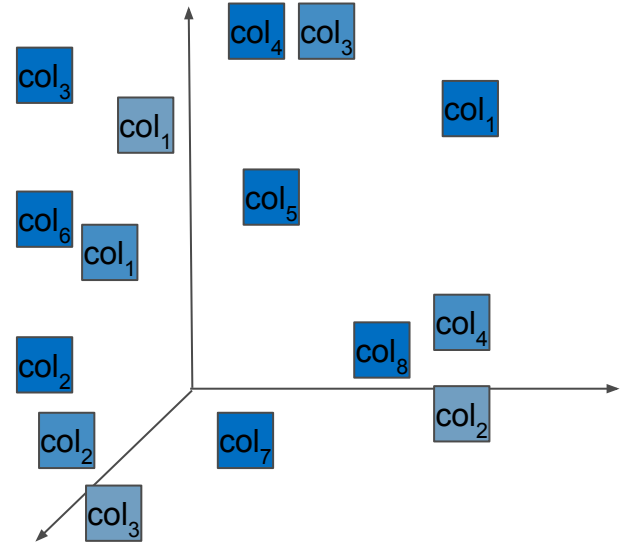
Pipeline to acquire data



Compressed Data Representation



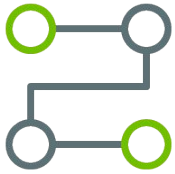
Context surrounding the data



Prolonged Value | Technology Solution



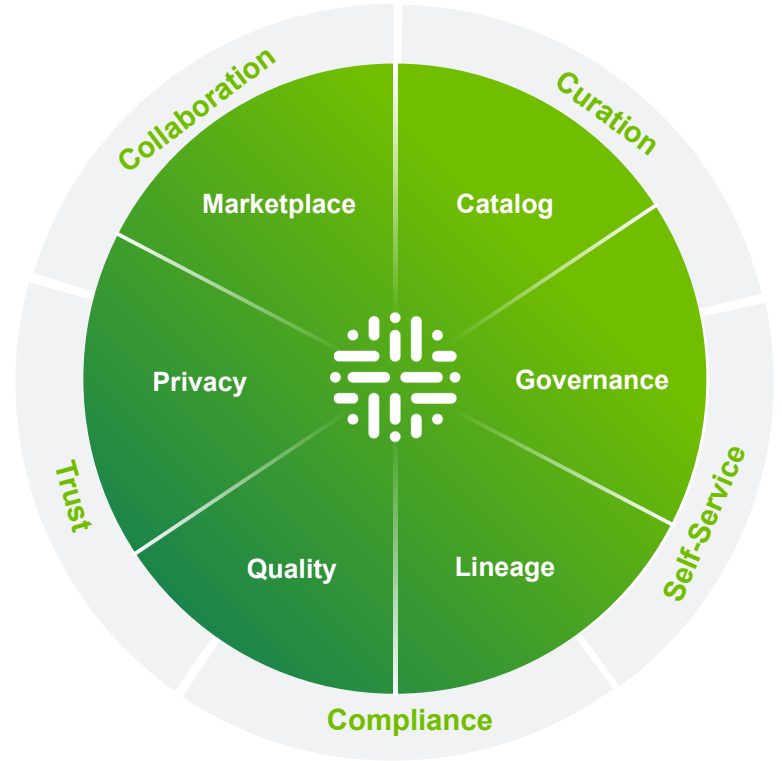
Pipeline to acquire data



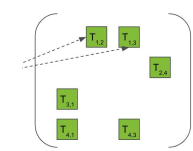
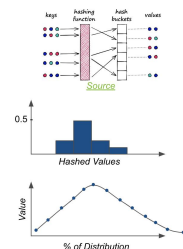
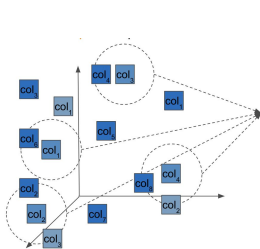
Compressed Data Representation



Context surrounding the data



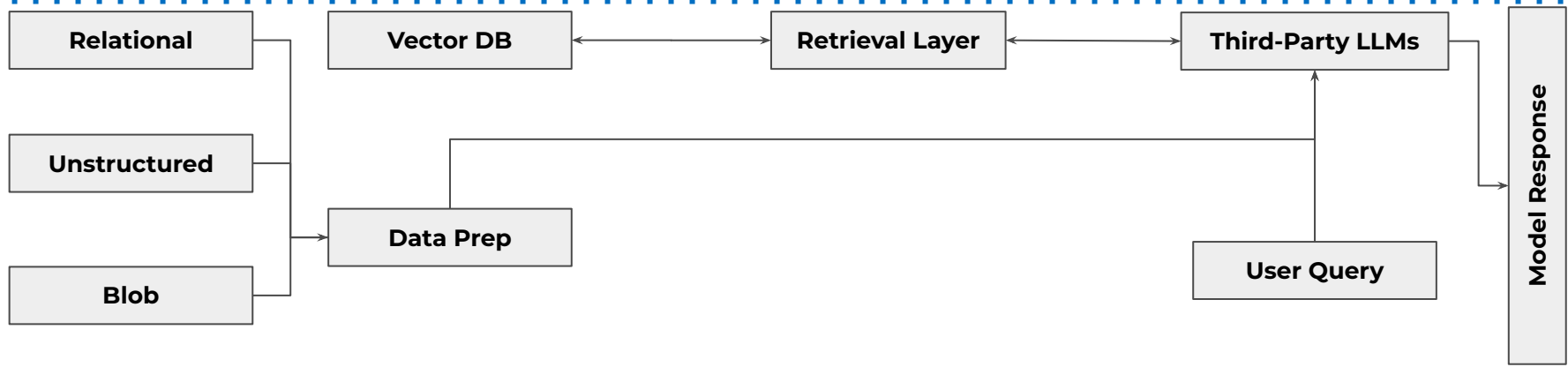
Expanding the Use-Case | What is RAG Anyways?



Data Layer

Compression Layer

Model Layer



Concluding Remarks

- Sometimes, you need a sledgehammer. Other times, you need a scalpel
- Finding the right people to navigate the various avenues of...
 - Mathematics
 - Statistics
 - ML
 - Neural Networks
 - LLMs...will drive pragmatic, performant solutions
- At Collibra, the AI team navigated those various bounds to deliver a data similarity solution which walks the tightrope of performance, compute, and scalability



[Imagine.art](#) 's interpretation of 'Data Council'

Thank you!