# Using Scalable Inference to Build DeepDive: A Declarative Dark Data System.

**Christopher Ré**

Stanford University

Computer Science

**Dark Data** System: ETL on Steroids

**Quality** that can exceed paid human annotators and volunteers

**DeepDive**

Extraction, Integration, & Cleaning
are *inference problems*

Focus on what matters:
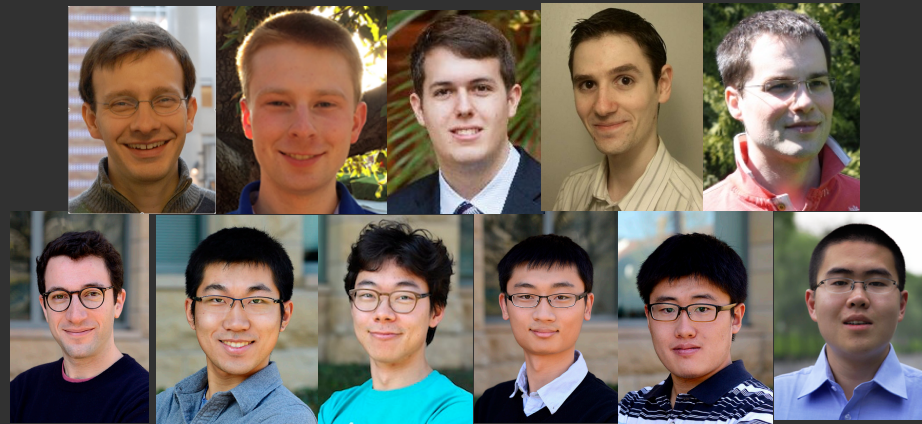*"Amdahl's law" for quality*

# DeepDive

Declarative Inference:
*Think about **features** not **algorithms**.*

Enables non-CS users, but
**scale** is a challenge.

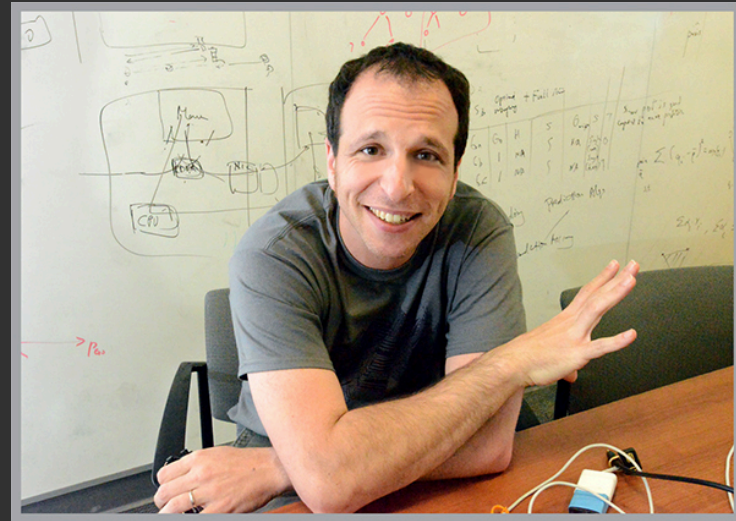# First, some **thank yous**

The DeepDive Team
http://deepdive.stanford.edu/

# INCOMPLETE THANK YOU

# Special Thank You to
# my optimization advisors

**Stephen J. Wright**
**Wisconsin.**
**God of Optimization.**

**Ben Recht**
**Berkeley.**
**Patriots Fan.**

# My Actual Advisor

# My Actual Boss

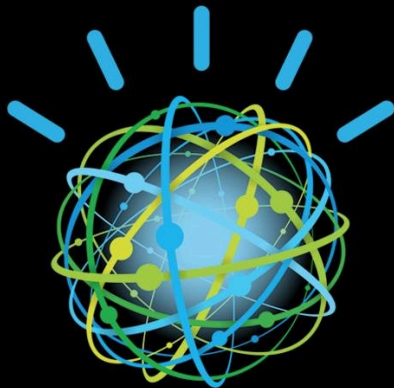# Many Existing Systems with the similar goal

Unstructured Input → Structured Knowledge Base

## Many Amazing Industrial Systems

IBM Research

Google

## An Incomplete Set of Awesome Research

TextRunner & ReVerb (Washington)

Knowledge Vault (Google)

liXto Lixto

YAGO-NAGA, SOFIE (MPI)

DBLife, xLog (Wisconsin)

StatSnowBall (Tshinghua & MSRA)

NELL (CMU)

ProbKB (Florida)

SystemT (IBM)

Many more

Back to our
**regularly scheduled programming…**

The world's scientific knowledge is **accessible**, but not **readable**.

# Today, some pressing problems require **macroscopic knowledge**
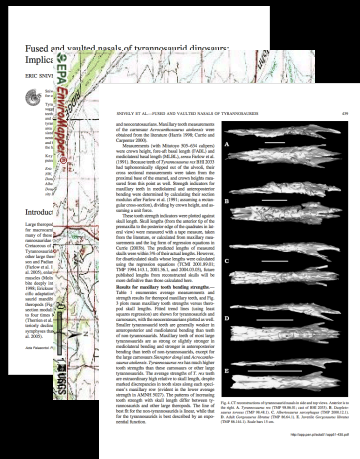


Climate & Biodiversity



Health



Financial Markets

# Could we build a machine to **read** for us?
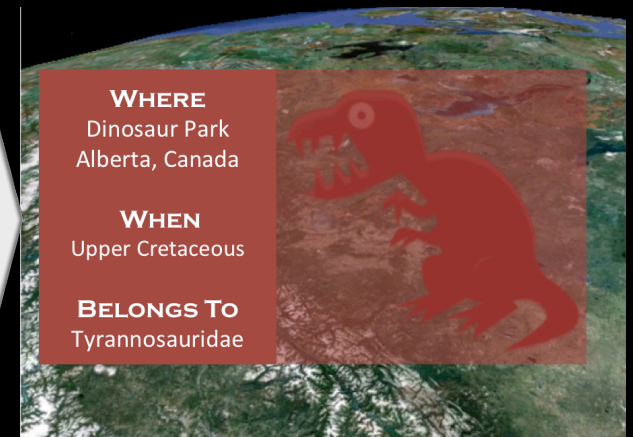
# PaleoDeepDive

## The Goal

Extract paleobiological facts to build **high coverage** fossil record.

T. Rex are found dating to the upper Cretaceous.

Statistical Inference
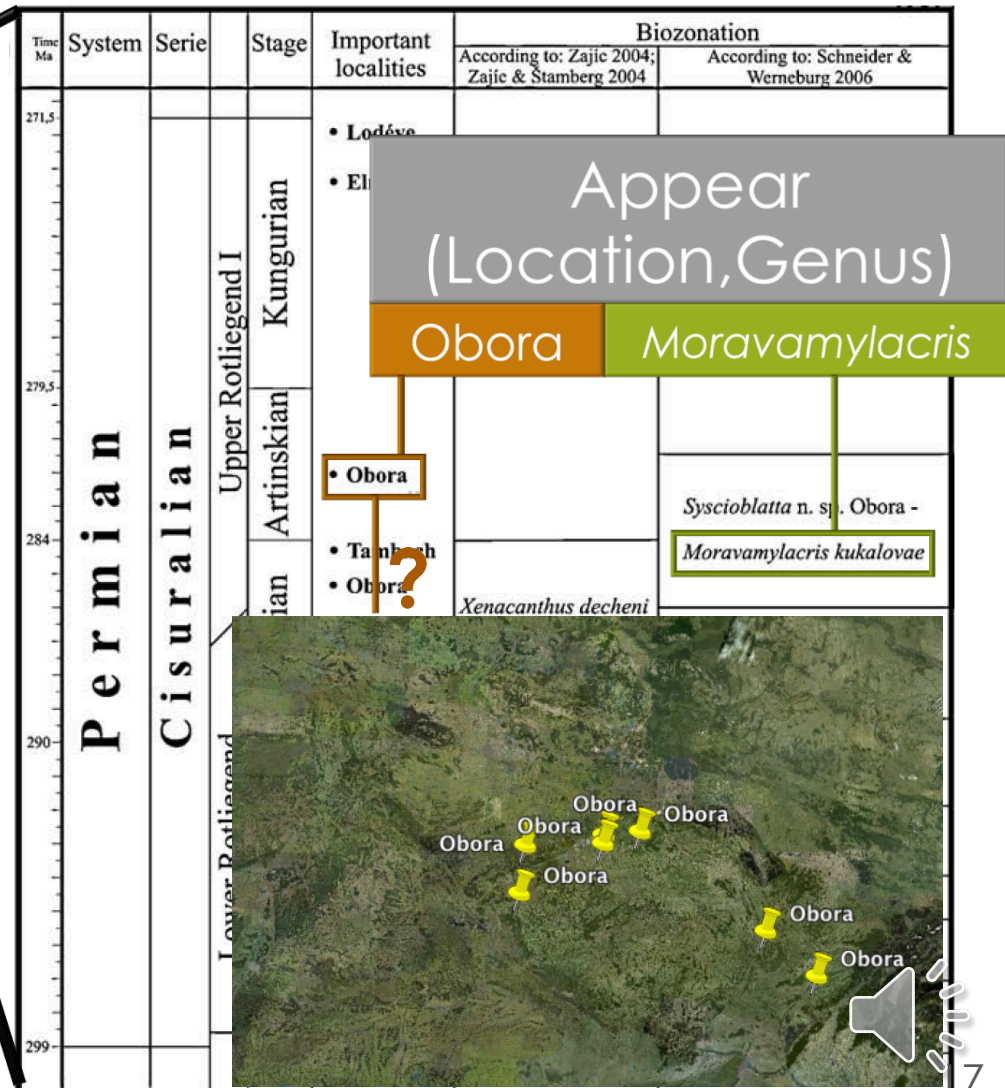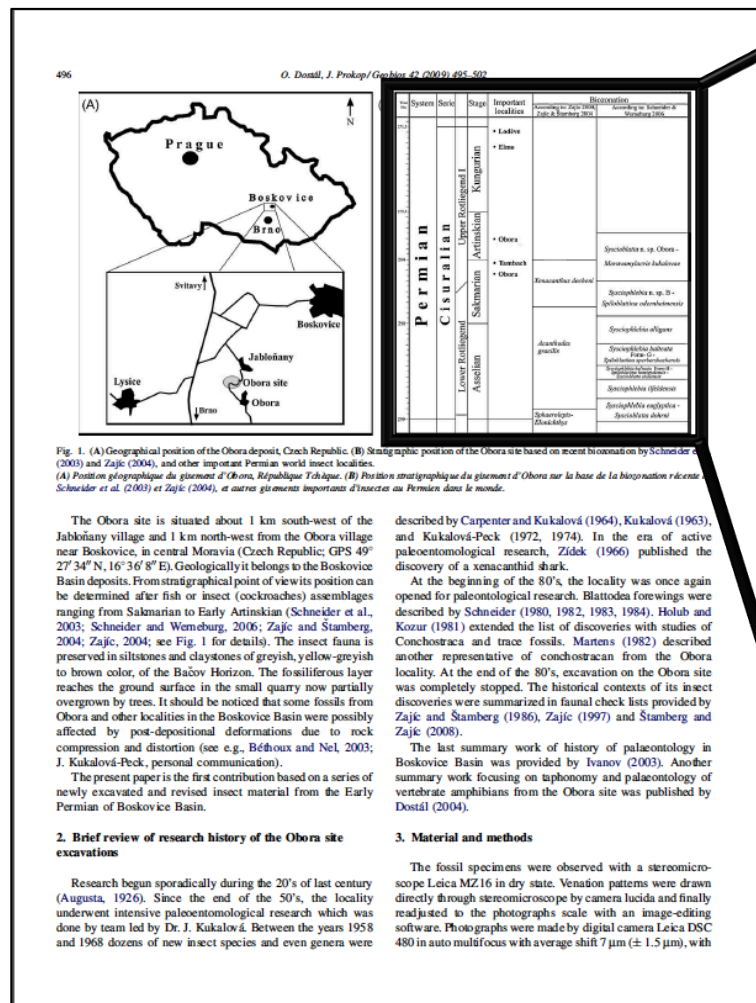
Appears("T. Rex", "Cretaceous")

**WHERE**
Dinosaur Park
Alberta, Canada

**WHEN**
Upper Cretaceous

**BELONGS TO**
Tyrannosauridae

## Aggressive Approach

*Every character, word, part of speech is a variable*
***Statistical inference*** *on billions of variables.*

16

Appear (Location, Genus)

Obora    Moravamylacris

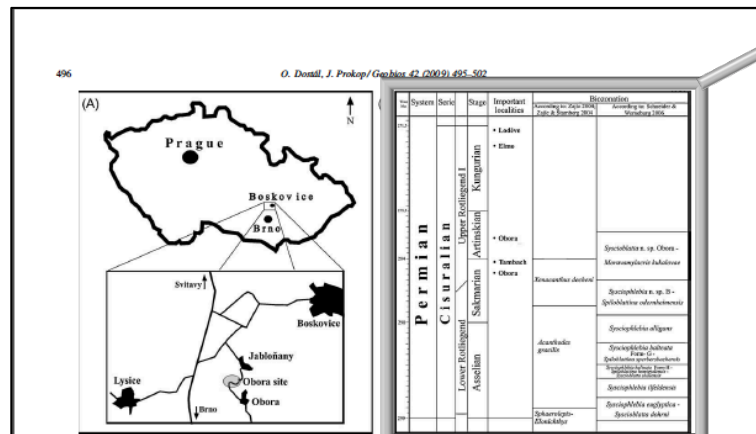# Data are buried in tables, but not in a self-contained way



Fig. 1. (A) Geographical position of the Obora deposit, Czech Republic. (B) Stratigraphic position of the Obora site based on recent biozonation by Schneider et al. (2003) and Zajíc (2004), and other important Permian world insect localities.
(A) Position géographique du gisement d'Obora, République Tchèque. (B) Position stratigraphique du gisement d'Obora sur la base de la biozonation récente de Schneider et al. (2003) et Zajíc (2004), et autres gisements importants d'insectes au Permien dans le monde.

The Obora site is situated about 1 km south-west of the Jabloňany village and 1 km north-west from the Obora village near Boskovice, in central Moravia (Czech Republic; GPS 49° 27′ 34″ N, 16° 36′ 8″ E). Geologically it belongs to the Boskovice

**Appear (Location, Genus)**

Obora — *Moravamylacris*

18

# **Joint** Probabilistic Inference Matters.

# PaleoDeepDive

Shanan Peters (Geo) and Miron Livny (CS)
DeepDive.Stanford.edu (**Ce Zhang** et al.)

# PaleoDB

# PaleoDeepDive

## Human-created

## Machine-created

329 volunteers
13 years
46K documents

10x **documents**.
100x **extractions**.

200+ Papers,
17 Nature/Science

Preliminary Precision

## Formation Precision

PaleoDB Volunteers: **0.84**

PaleoDeepDive: **0.94**

Peters, S., Zhang, C, Livny, M., and Ré, C. A New Machine-Aggregated Empirical History of Life on Earth. *PLOS ONE*, 2014 featured in *Nature* July 1, 15

# Hope: knowledge bases can help **accelerate science**.

Tree of Life      Drug Repurposing      Genomics

*Used by a number of companies with quality that best **professional** human annotators; winner of TACKBP14.*

# Human Trafficking on the (Dark) Web...

**DARPA MEMEX**

**Hypothesis**: Trafficked individuals offer *lower cost* and *riskier* sexual services.

*In Plain sight*: Web ads for such services

**Challenges**:
1. Need **high-resolution information** to build model.
   - *services for what rate, ethnicity, location, etc.*
2. Scientific papers are **clear**—*dark web is* **obfuscated**.

# Human Trafficking on the (Dark) Web…

# In Use by Law Enforcement

*New York DA use MEMEX Data for all trafficking investigations this year.* **Real Arrests**

For DARPA MEMEX, we were operational in 6 months

- Processed >35M documents (~26M records)
- Tens of columns (location, phone #, price, etc)
- With compute times of less than a day
- >90% Precision for most relations
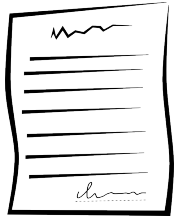
Declarative language allows **algorithmic independence**

# Example: Extracting Spouse Relations

**Corpus (Dark Data)**

U.S. President **Barack Obama**'s wife **Michelle Obama** honored all mothers on Mother's Day and offered her thoughts …

## How do we produce tuples like

`Married(Barack Obama, Michelle Obama)`

## **And** all other married couples in text?

*Examples on Github!*

# 0. Data Preprocessing

**Corpus (Dark Data)**

**text**

U.S. President Barack Obama's wife Michelle Obama honored all mothers on Mother's Day and offered her thoughts …

**DDlog**: **Declarative** inspired by Datalog, MLNs

```
Sentences :- !nlp(Corpus).
function nlp over (text)
        returns (words,pos,ner,sid)
  implementation "udf/corenlp.sh …" .
```

**Sentences**

| words | POS | NER | SID |
|---|---|---|---|
| [U.S.,President,Barack,Obama,'s,wife,Michelle,Obama ,…] | [NNP,NNP,NNP,NNP, POS,NN,NNP,NNP, …] | [LOC,O,PER,PER,O,O,PER,PER, …] | S1 |

**Sentences**

| words | POS | NER | SID |
|---|---|---|---|
| [U.S.,President,Barack,Obama,'s,wife,Michelle,Obama,...] | [NNP,NNP,NNP,NNP,POS,NN,NNP,NNP,...] | [LOC,O,PER,PER,O,O,PER,PER,...] | S1 |

```
Mentions :- !ext_person(Sentences).
function ext_person over (words,pos,ner,sid)
                returns (sid,mid,words)
        implementation "udf/find_person.py".
```

```
MarriedCandidate(s,p1,p2) :-
   Mentions(s,p1,_), Mentions(s,p2,_).
```

**Mentions**

| SID | MID | words |
|---|---|---|
| S1 | M1 | [Barack,Obama] |
| S1 | M2 | [Michelle,Obama] |

**MarriedCandidate**

| SID | MID | MID |
|---|---|---|
| S1 | M1 | M2 |

29

# 2. Feature Extraction

**Sentences**

| words | POS | NER | SID |
|---|---|---|---|
| [U.S.,President,Barack,Oba ma,'s,wife,Michelle,Obama ,...] | [NNP,NNP,NNP,NNP, POS,NN,NNP,NNP, ...] | [LOC,O,PER,PE R,O,O,PER,PER, ...] | S1 |

```
Features :- !ext_features(Sentences, MarriedCandidate).
function ext_features …
    implementation "udf/ext_features.py".
```

**Mentions**

| SID | MID | words |
|---|---|---|
| S1 | M1 | [Barack,Obama] |
| S1 | M2 | [Michelle,Obama] |

**MarriedCandidate**

| SID | MID | MID |
|---|---|---|
| S1 | M1 | M2 |

**Features**

| MID | MID | feature |
|---|---|---|
| M1 | M2 | 's wife |

**Sentences**

| words | POS | NER | SID |
|---|---|---|---|
| [U.S.,President,Barack,Obama,'s,wife,Michelle,Obama,...] | [NNP,NNP,NNP,NNP,POS,NN,NNP,NNP,...] | [LOC,O,PER,PER,O,O,PER,PER,...] | S1 |

```
Married(p1,p2) :-
    MarriedCandidate(_,p1,p2),
    Features(p1,p2,f)
weight = f.
```

Just defined a Binary classifier!

*Married is an (incomplete) set of examples*

**Mentions**

| SID | MID | words |
|---|---|---|
| S1 | M1 | [Barack,Obama] |
| S1 | M2 | [Michelle,Obama] |

**MarriedCandidate**

| SID | MID | MID |
|---|---|---|
| S1 | M1 | M2 |

**Features**

| MID | MID | feature |
|---|---|---|
| M1 | M2 | 's wife |

31

# Users write Features & Transformation in **DDlog** (Inspired by MLNs) & Python.

```
Sentences :- !nlp(Corpus).
function nlp over (text) returns (words,pos,ner,sid)
    implementation "udf/corenlp.sh …" .


Mentions :- !ext_person(Sentences).
function ext_person over (words,pos,ner,sid) returns (sid,mid,words)
            implementation "udf/find_person.py".
MarriedCandidate(s,p1,p2) :- Mentions(s,p1,_), Mentions(s,p2,_).


Features :- !ext_features(Sentences, MarriedCandidate).
function ext_features … implementation "udf/ext_features.py".


Married(p1,p2) :- MarriedCandidate(_,p1,p2), Features(p1,p2,f)
weight = f.
```
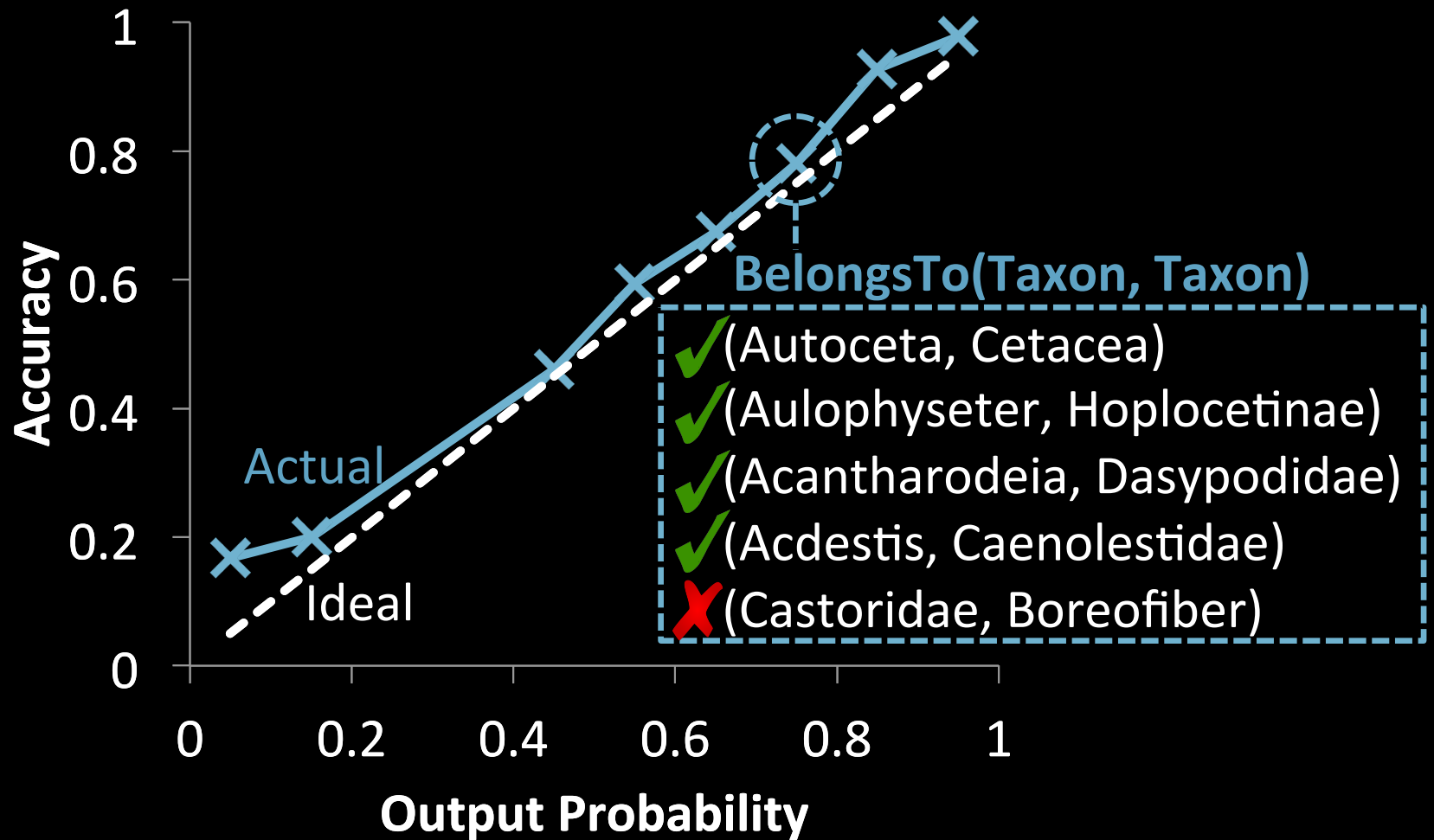
**No** reference to algorithms: Just features.
DD does the rest. *(Demos online)*

# Algorithmic Independence

Define the *meaning* **independently** of *the algorithm used to compute it.*

**algorithmic independences,**
*requires a fast* engine…..

**Key Issue:** Balance _Statistical_ versus _Hardware_ Efficiency.

*Ce Zhang* *DimmWitted: A Study of Main-Memory Statistical Analytics. VLDB14.*

# Statistical Analytics Crash Course

Staggering amount of machine learning/stats can be written as:

$$\min_x \sum_{i=1}^{N} f(x, y_i)$$

N (number of $y_i$s, data) typically in the billions
Ex: Classification, Recommendation, Deep Learning.

*De facto* iteration to solve large-scale problems: **SGD**.

$$x^{k+1} = x^k - \alpha N \nabla f(x^k, y_j)$$

Select one term, j, and estimate gradient.

Billions of tiny iterations.

# How do we run SGD in Parallel?

**Data Systems Perspective of SGD.**

$$x^{k+1} = x^k - \alpha N \nabla f(x^k, y_j)$$

**Insane conflicts:** Billions of tiny (~100 instructions) jobs, RW conflicts on $x$, which is called **the model.**

How can we hope to speed this up with parallelism?

Serializability seems hopeless...

**Thm:** If we do **no locking**, SGD still converges to right answer—at essentially the same theoretical rate!

**Hogwild**! [Niu, Recht, **Ré,** Wright NIPS11]
**AsySCD** [Liu, Wright et al. ICML14, JMLR14]
**Buckwild! [**DeSa et al. ICML15**]**

Technical conditions on ratio of processors, delays, (semantic) sparsity.

**High-level idea**: Go Hogwild! answer is only *statistically* correct.

# A larger trend?

*NB: There is theory here SGD [NIPS11,NIPS12], SCD [ICML14,ICML15], more soon and systems work [SIGMOD13, SIGMOD14, VLDB14]

Relaxing **consistency** to be **architecturally aware** can be a big performance win.
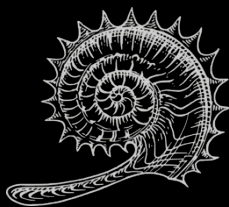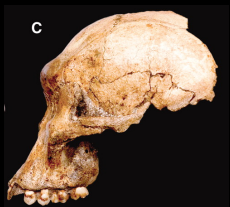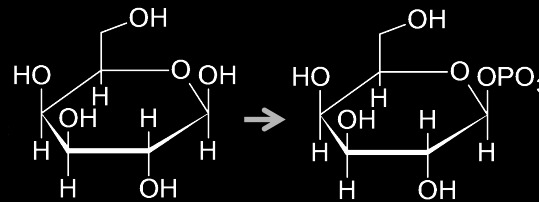
# Dark(er) Data Systems

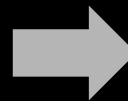# *Everyone*: Broadly Usable

**PaleoDeepDive**
- 🕐 2 years
- 👤 1 CS Student

→

**PharmaDeepDive**
- 🕐 6 months
- 👤 1 BioE Student

How do we make building a KB **easier and cheaper**?

Think about **features**, not **algorithms**.

A **framework** for feature engineering.
*[SIGMOD14: train 100 models as quickly as 1]*

# Conclusion

1. **Dark Data** to help with **macroscopic questions**

2. Probabilistic inference = **algorithmic** independence

3. **Hardware v. Statistical** Efficiency.