# From Data Fusion to Knowledge Fusion

*Xin Luna Dong,* Evgeniy Gabrilovich, Geremy Heitz, Wilko Horn,
Kevin Murphy, Shaohua Sun, Wei Zhang

News

# Google's Knowledge Vault already contains 1.6 billion facts
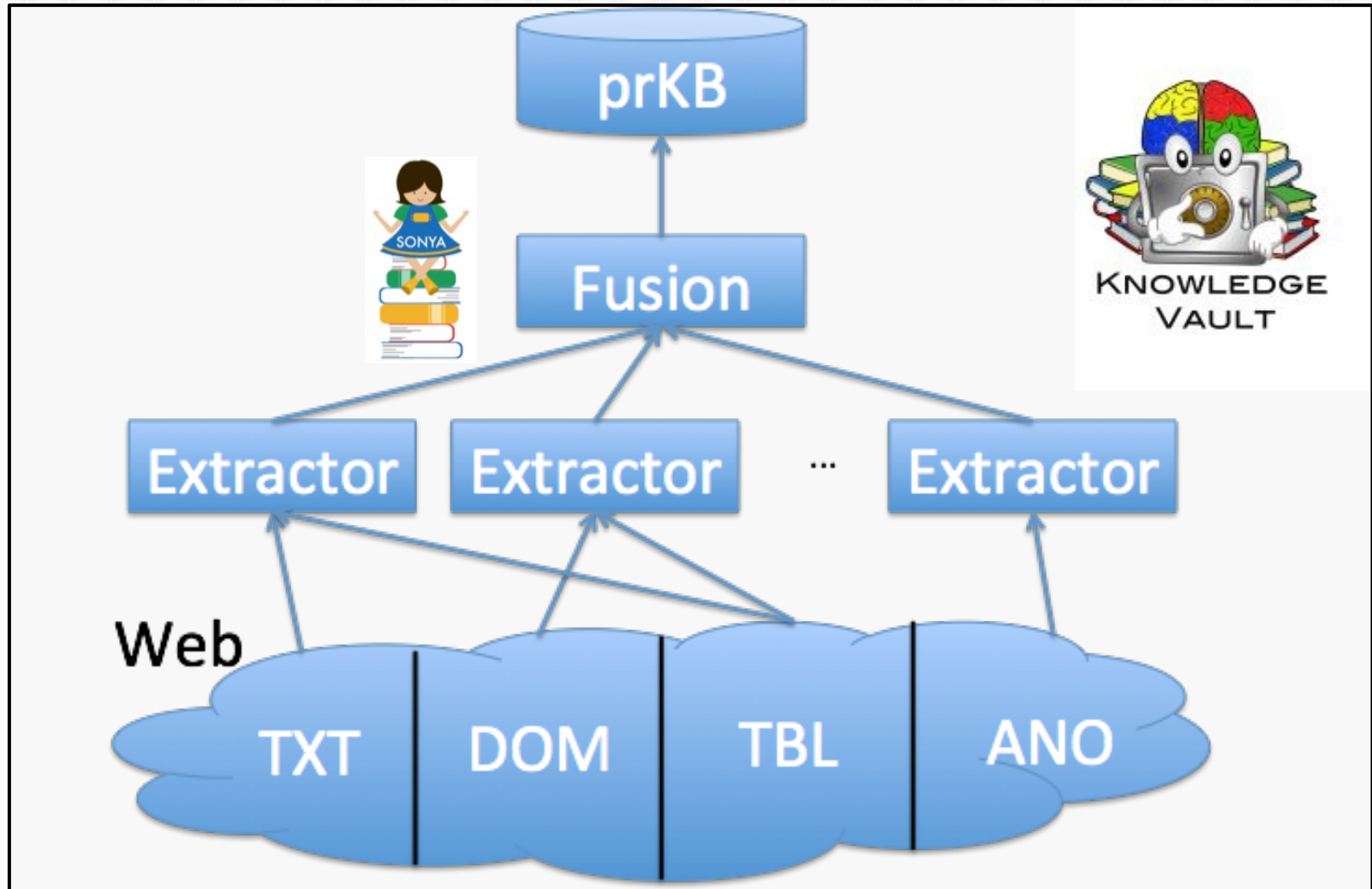
FELICITY NELSON
SATURDAY, 23 AUGUST 2014

The automated, fact-harvesting bot will build up a collection of all human knowledge.

# Knowledge Vault– Building a Probabilistic KB

[**VLDB'2014**, Sigmod'2014, KDD'2014]

# Extracted Knowledge

- Triple: (subject, predicate, object)
  e.g., (Tom Cruise, date_of_birth, 7/3/1962)
  - Subject–a Freebase mid
    e.g., /m/07r1h
  - Predicate–predefined in
    Freebase; e.g., people/person/date_of_birth
  - Object–a Freebase mid, a string, a number, or a date.

# Observation I.

Yes, We CAN Build A Large KB from the Web! :-)

# Statistics for Extracted Triples

● A large knowledge base

As of 11/2013

| #Triples | 1.6B (now 2.8B) |
|---|---|
| #Subjects (Entities) | 43M |
| #Predicates | 4.5K |
| #Objects | 102M |

# Observation II.
# But, A Lot of Mistakes :-(

# Errors Can Creep in at Every Stage

Extraction error: (Obama, nationality, Chicago)

# Errors Can Creep in at Every Stage

Reconciliation error:
(Obama, nationality, North America)



American President Barack Obama

9/2013

# Errors Can Creep in at Every Stage

Source data error: (Obama, nationality, Kenya)



Obama born in Kenya

# Statistics for Triple Correctness

- The gold standard (based on Freebase) contains about 40% of the triples
- Overall accuracy: 30%
- Random sample on 25 false triples
  - Extraction errors: 24 (96%)
  - Source-data errors: 1 (4%)

# Knowledge Fusion

- Input: Knowledge triples and their provenances (i.e., which extractor extracts from which source)
- Output: a probability in [0,1] for each triple
  - High pr→ search, etc.
  - Medium pr→ active learning, probabilistic inference, etc.
  - Low pr→ Negative training examples

# Observation III.
*Data Fusion* Techniques Work Fairly Well for *Knowledge Fusion*

# Data Fusion–Definition

**Input**

Sources

| | $S_1$ | $S_2$ | ... | $S_N$ |
|---|---|---|---|---|
| $D_1$ | | | | |
| $D_2$ | | | | |
| $D_3$ | | | | |
| ... | | | | |
| $D_M$ | | | | |

Data items

**Output**

Truths

| | |
|---|---|
| $D_1$ | |
| $D_2$ | |
| $D_3$ | |
| ... | |
| $D_M$ | |

Data items

# Data Fusion–Intuition

|  | Src1 | Src2 | Src3 |
|---|---|---|---|
| Jagadish | UM | ATT | UM |
| Dewitt | MSR | MSR | UW |
| Bernstein | MSR | MSR | MSR |
| Carey | UCI | ATT | BEA |
| Franklin | UCB | UCB | UMD |

# Data Fusion–Intuition

| | Src1 | Src2 | Src3 |
|---|---|---|---|
| Jagadish | UM | ATT | UM |
| Dewitt | MSR | MSR | UW |
| Bernstein | MSR | MSR | MSR |
| Carey | UCI | ATT | BEA |
| Franklin | UCB | UCB | UMD |

Voting--Trust the majority.

# Data Fusion–Intuition

|  | Src1 | Src2 | Src3 |
|---|---|---|---|
| Jagadish | UM | ATT | UM |
| Dewitt | MSR | MSR | UW |
| Bernstein | MSR | MSR | MSR |
| Carey | UCI | ATT | BEA |
| Franklin | UCB | UCB | UMD |

# Data Fusion–Intuition

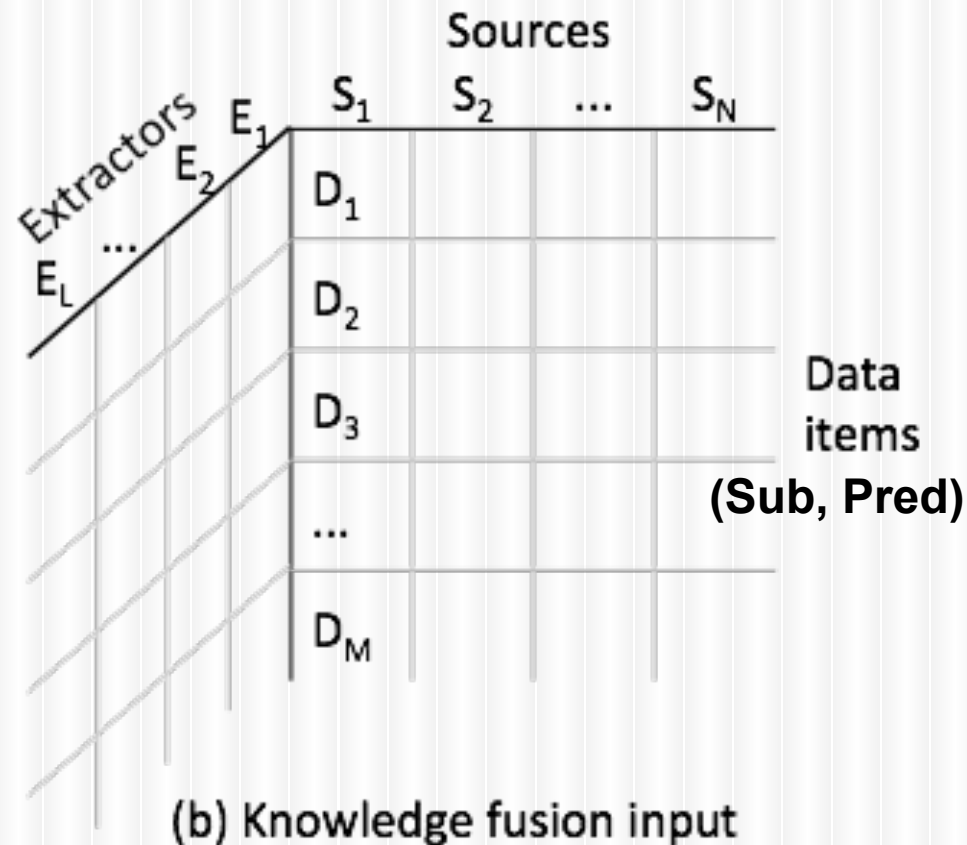|           | Src1 | Src2 | Src3 |
|-----------|------|------|------|
| Jagadish  | UM   | ATT  | UM   |
| Dewitt    | MSR  | MSR  | UW   |
| Bernstein | MSR  | MSR  | MSR  |
| Carey     | UCI  | ATT  | BEA  |
| Franklin  | UCB  | UCB  | UMD  |

Quality-based--Give higher votes to more accurate sources.

## I. Input is *three-dimensional*



(a) Data fusion input

(b) Knowledge fusion input
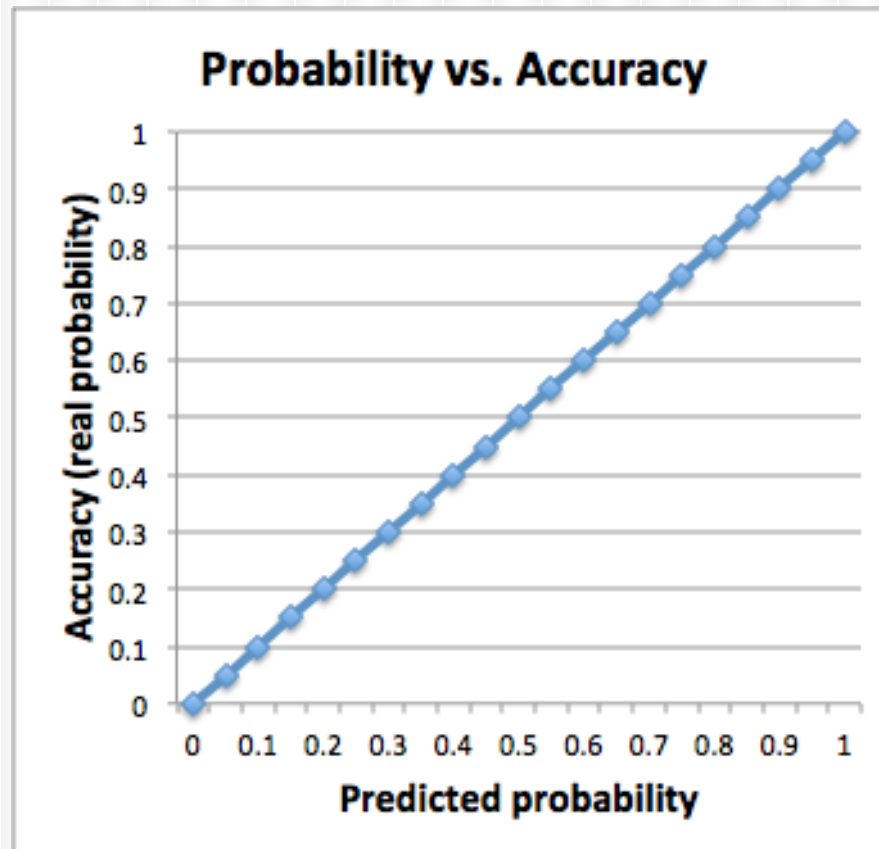
## II. Output prs should be *well-calibrated*



Probability vs. Accuracy — Accuracy (real probability) vs. Predicted probability

# Knowledge Fusion Challenges

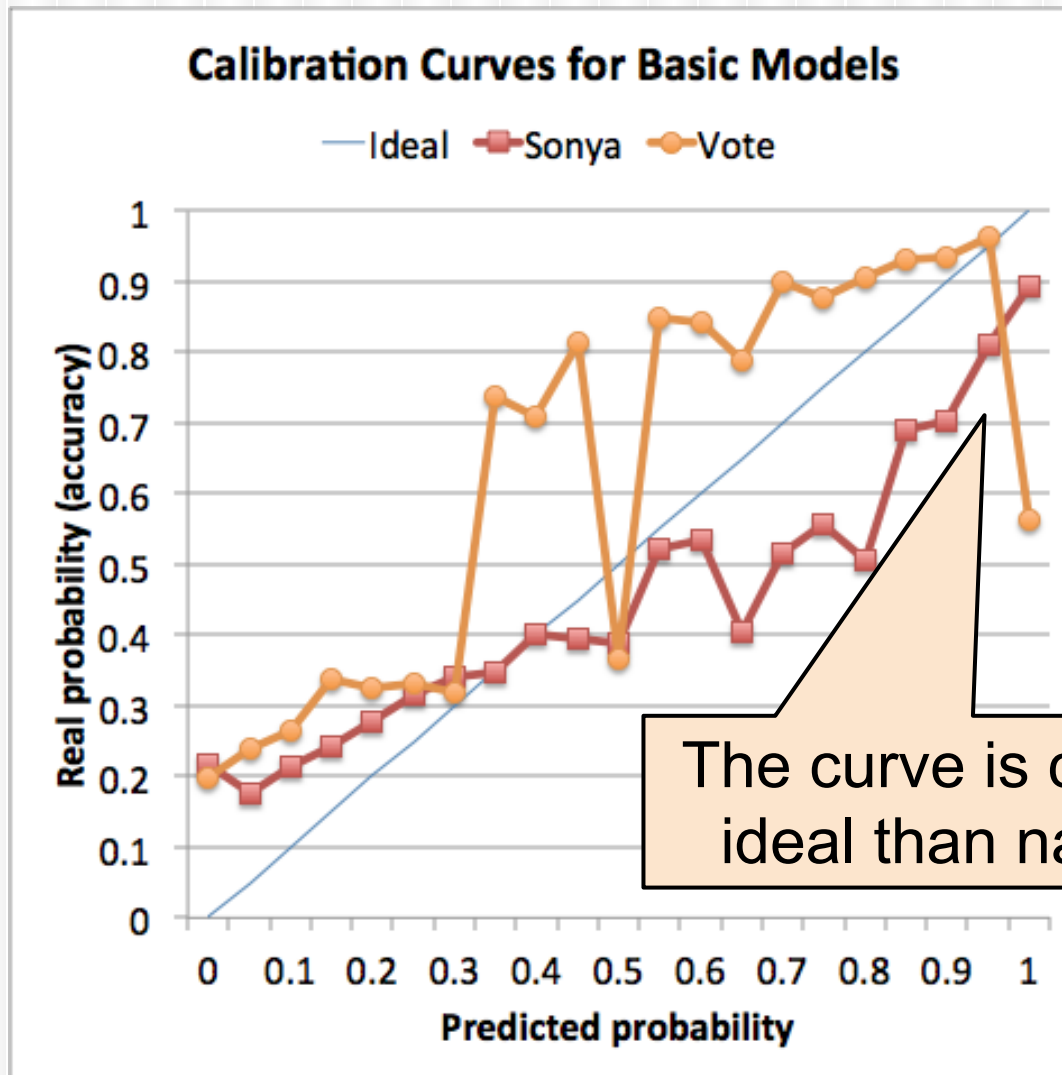## III. Data are of *Web-scale*

- Three orders of magnitude larger than currently published data-fusion applications
  - Size: 1.1TB
  - Sources: 170K→ 1B+
  - Data items: 400K→375M
  - Values: 18M→6.4B (1.6B unique)
- Data are highly skewed
  - #Triples/Data-item: 1 - 2.7M
  - #Triples/Source: 1 - 50K

# Knowledge Fusion Solutions

- Treat each (URL, Extractor) as a whole *(provenance)* for accuracy evaluation
- A series of refinements to a Bayesian model to improve probability calibration
- MapReduce Based Framework
  - Sample for *too big* data items or provenances

# Basic Sonya Solution vs. Voting



**Calibration Curves for Basic Models**

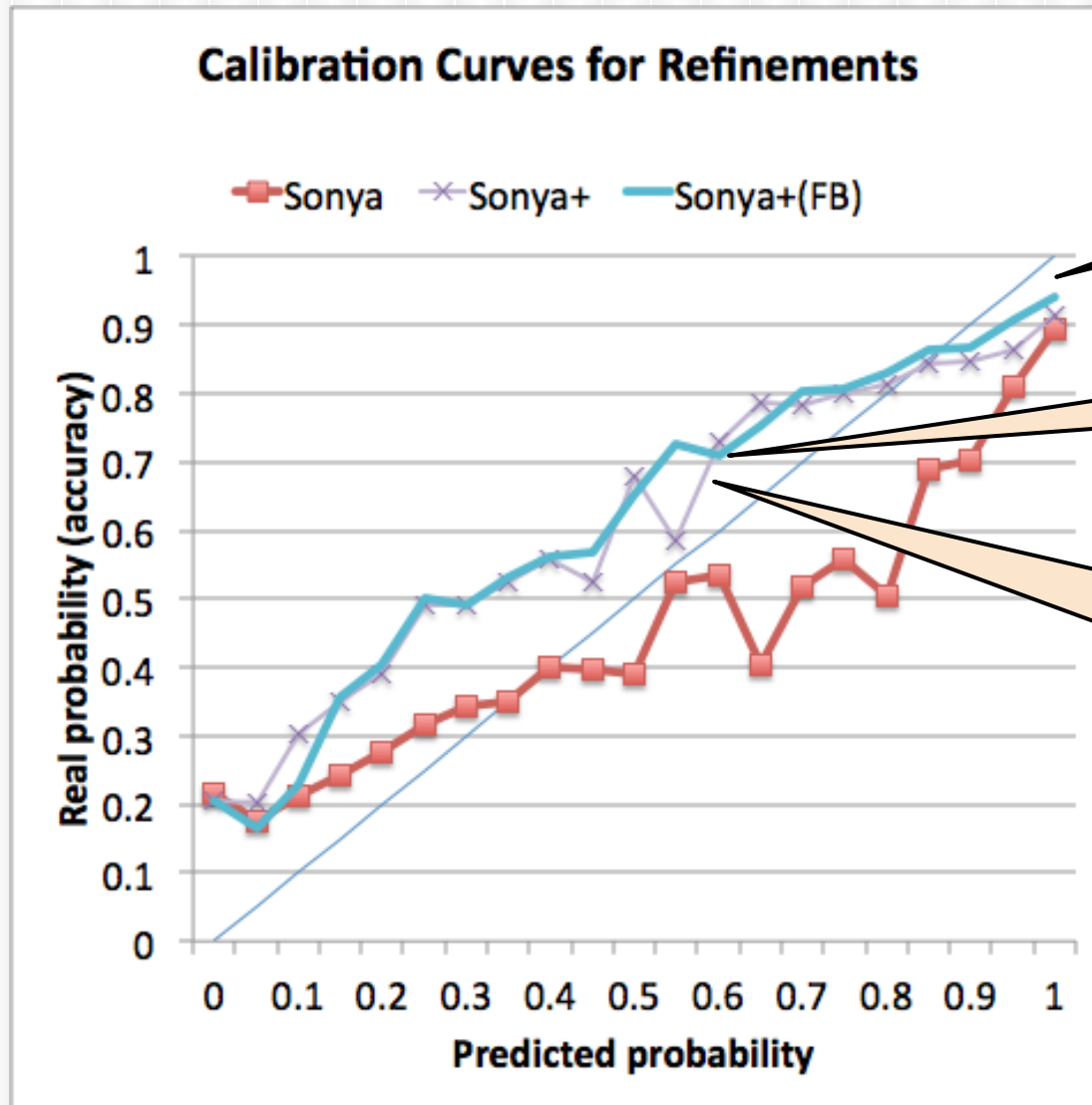The curve is closer to the ideal than naive voting

# Refinements

I. Ignore low-coverage provenances
II. Granularity (URL->Site, Extractor->Pattern, Predicate)
III. Ignore low-accuracy provenances
IV. Initiate provenance accuracy by FB

+I, II, III. Sonya+ : unsupervised
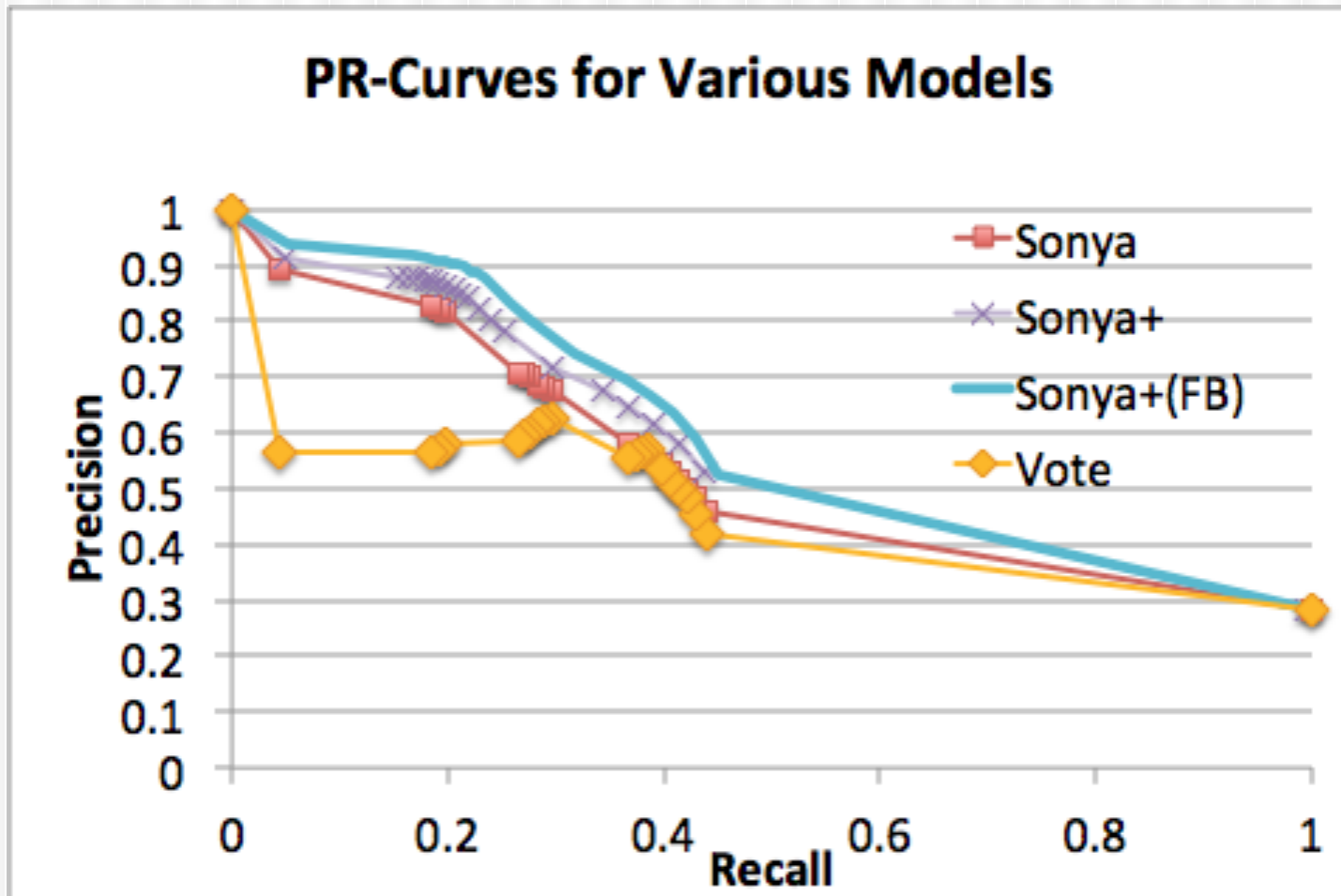+IV. Sonya+(FB) : semi-supervised

# Calibration Curve



Calibration Curves for Refinements

Legend: Sonya, Sonya+, Sonya+(FB)

Y-axis: Real probability (accuracy)
X-axis: Predicted probability

Higher accuracy

Smoother curve

Unsupervised methods very effective

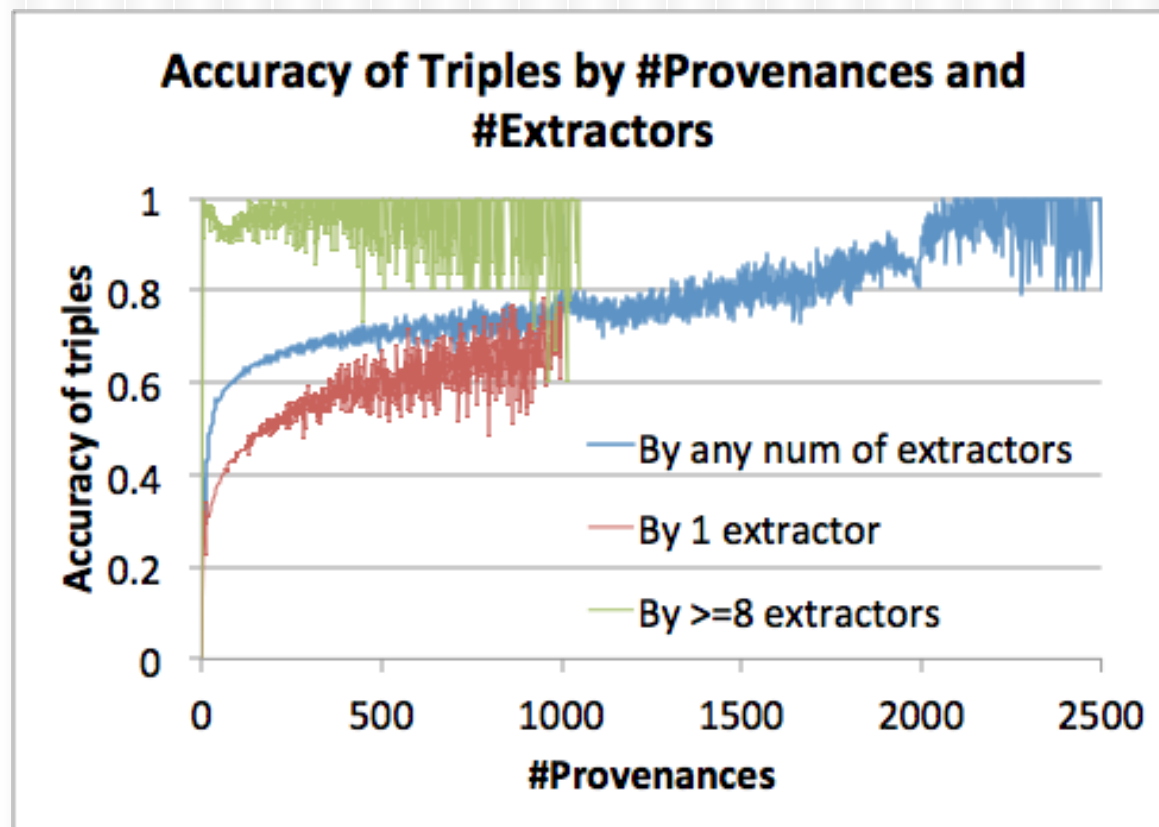Weighted deviation: $0.037 \rightarrow 0.035 \rightarrow 0.032$

# Precision-Recall Curve

# Observation IV.
# Still Many Places to Improve!

# One Inherent Limitation

Cannot distinguish errors from extractors and from sources



Accuracy of Triples by #Provenances and #Extractors

# TAKE AWAYS

- A new area--Knowledge Fusion

- We can solve KF problem fairly well by adapting DF methods

- Many interesting future directions for KF!

- Many exciting applications for the prKB!!

# THANK YOU!

*Questions?*