

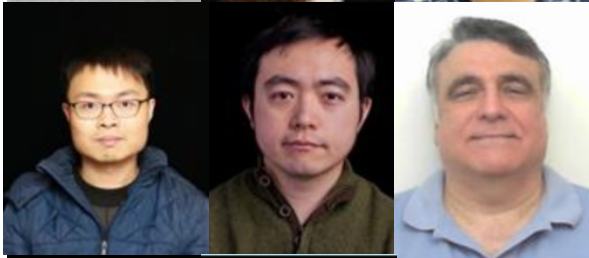
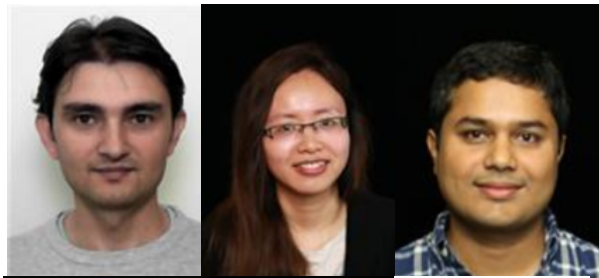
Data Integration and Machine Learning: A Natural Synergy

Xin Luna Dong @ Amazon.com

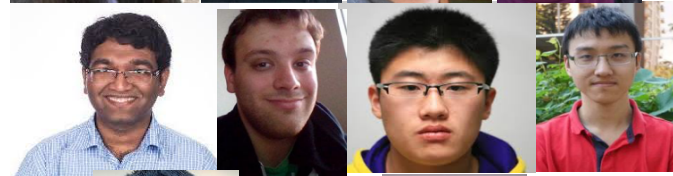
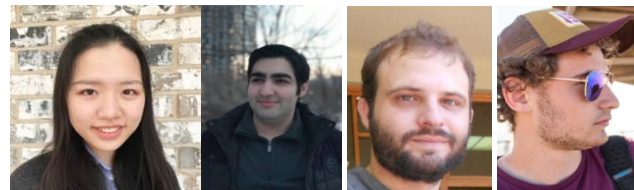
Theo Rekatsinas @ UW-Madison

Sigmod 2018

Acknowledgement



snorkel



What is Data Integration?

- **Data integration:** to provide unified access to data residing in multiple, autonomous data sources
 - **Data warehouse:** create a single store (materialized view) of data from different sources offline. Multi-billion dollar business.
 - **Virtual integration:** support query over a mediated schema by applying online query reformulation. E.g., Kayak.com.
- In the RDF world: different names for similar concepts
 - **Knowledge graph** is equivalent to a data warehouse. Has been widely used in Search and Voice
 - **Linked data** is equivalent to virtual integration

Why is Data Integration Hard?

- Heterogeneity everywhere
 - Different data formats

Web tables & Lists

	Name and (party) ¹
1.	Washington
2.	J. Adams
3.	Jefferson
4.	Madison

DOM Trees

Free texts

Diagram

	Biological level	Examples	Pre-amputation	Post-amputation	Regenerate
Whole body		Regeneration from a small body fragment			
Structure		Limb, fin, tail, head, tentacle, siphon, arm, stalk			
Internal organ		Heart, liver, lens			
Tissue		Epidermis, gut lining			
Cell		Axon, muscle fiber			

Data Extraction



Schema Alignment



Entity Linkage



Data Fusion

Why is Data Integration Hard?

- Heterogeneity everywhere
 - Different ways to express the same classes and attributes

IMDB



Anahí

[Actress](#) | [Music Department](#) | [Soundtrack](#)

Anahi was born in Mexico. She's had roles in Tu y Yo, in which she played a 17 year old girl while she was 13, and Vivo Por Elena, in which she played Talita, a naive and innocent teenager. Anahi lives with her mother and sister name Marychelo. She hopes to become a fashion designer one day, and is currently pursuing a career in singing.

[See full bio »](#)

Born: May 14, 1982 in Mexico City, Distrito Federal, Mexico

[More at IMDbPro »](#)

[Contact Info: View manager](#)



SEE RANK

WikiData

Anahí Puente (Q1694)

Mexican singer-songwriter and actress
Mia

[In more languages](#) [Configure](#)

Language	Label
English	Anahí Puente
Chinese	阿纳希·普恩特
Spanish	Anahí Puente

[date of birth](#)

7 November 198

[1 reference imported from](#)

Data Extraction



Schema Alignment



Entity Linkage

No description defined

Cantante, compositora y actriz mexicana



Data Fusion

+ add value

Why is Data Integration Hard?

- Heterogeneity everywhere
 - Different references to the same entity

IMDB



Anahí [SEE RANK](#)

[Actress](#) | [Music Department](#) | [Soundtrack](#)

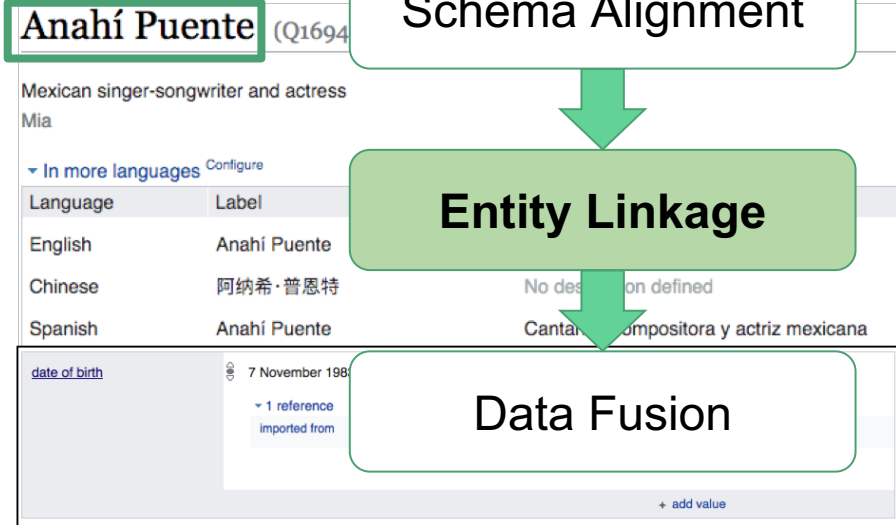
Anahi was born in Mexico. She's had roles in *Tu y Yo*, in which she played a 17 year old girl while she was 13, and *Vivo Por Elena*, in which she played Talita, a naive and innocent teenager. Anahi lives with her mother and sister name Marychelo. She hopes to become a fashion designer one day, and is currently pursuing a career in singing. [See full bio »](#)

Born: May 14, 1982 in Mexico City, Distrito Federal, Mexico

[More at IMDbPro »](#)

Contact Info: [View manager](#)

WikiData




Anahí Puente (Q1694)

Mexican singer-songwriter and actress
Mia

[In more languages](#) [Configure](#)

Language	Label
English	Anahí Puente
Chinese	阿纳希·普恩特
Spanish	Anahí Puente

[date of birth](#)  7 November 1981

[1 reference imported from](#)

[+ add value](#)

Data Extraction



Schema Alignment



Entity Linkage

No description defined

Cantante, compositora y actriz mexicana



Data Fusion

Why is Data Integration Hard?

- Heterogeneity everywhere
 - Conflicting values

IMDB



Anahí

[Actress](#) | [Music Department](#) | [Soundtrack](#)



Anahi was born in Mexico. She's had roles in Tu y Yo, in which she played a 17 year old girl while she was 13, and Vivo Por Elena, in which she played Talita, a naive and innocent teenager. Anahi lives with her mother and sister name Marychelo. She hopes to become a fashion designer one day, and is currently pursuing a career in singing.

[See full bio »](#)

Born: May 14, 1982 in Mexico City, Distrito Federal, Mexico

[More at IMDbPro »](#)

[Contact Info: View manager](#)

WikiData

Anahí Puente (Q1694)

Mexican singer-songwriter and actress
Mia

[In more languages](#) [Configure](#)

Language	Label
English	Anahí Puente
Chinese	阿纳希·普恩特
Spanish	Anahí Puente

[date of birth](#)

7 November 1982

[1 reference imported from](#)

[+ add value](#)

Data Extraction



Schema Alignment



Entity Linkage

No description defined

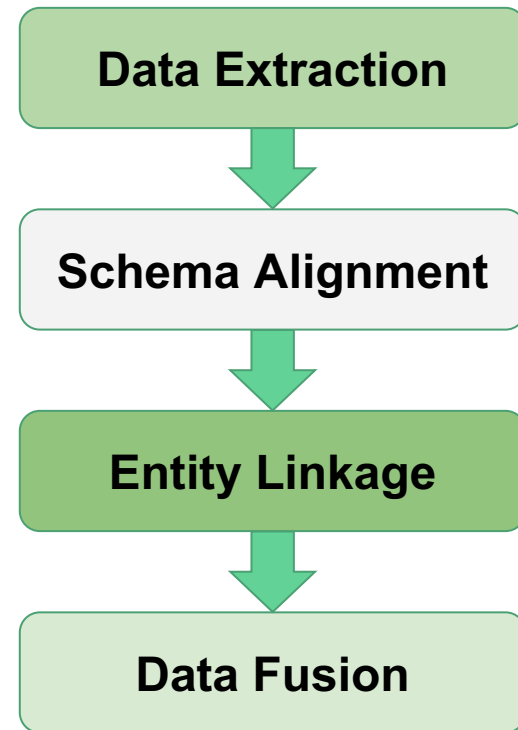
Cantante, compositora y actriz mexicana



Data Fusion

Importance from a Practitioner's Point of View

- Entity linkage is indispensable whenever integrating data from different sources
- Data extraction is important for integrating non-relational data
- Data fusion is necessary in presence of erroneous data
- Schema alignment is helpful when integrating relational data, but not affordable for manual work if we integrate many sources



What is Machine Learning?

- **Machine learning:** teach computers to *learn* with data, not by programming

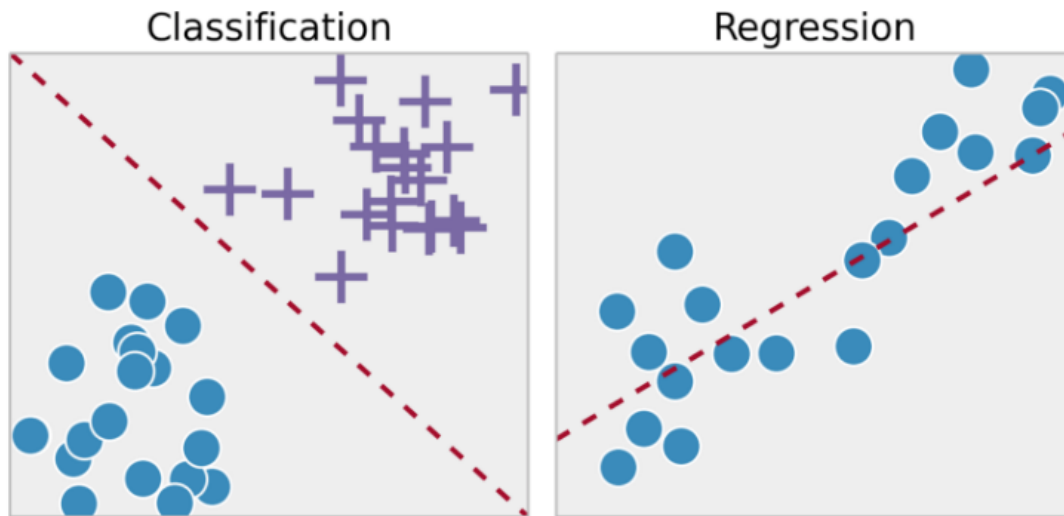
- **More Formal definition**

A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P if its performance at tasks in T , as measured by P , **improves with experience E .**

-- Tom Mitchell

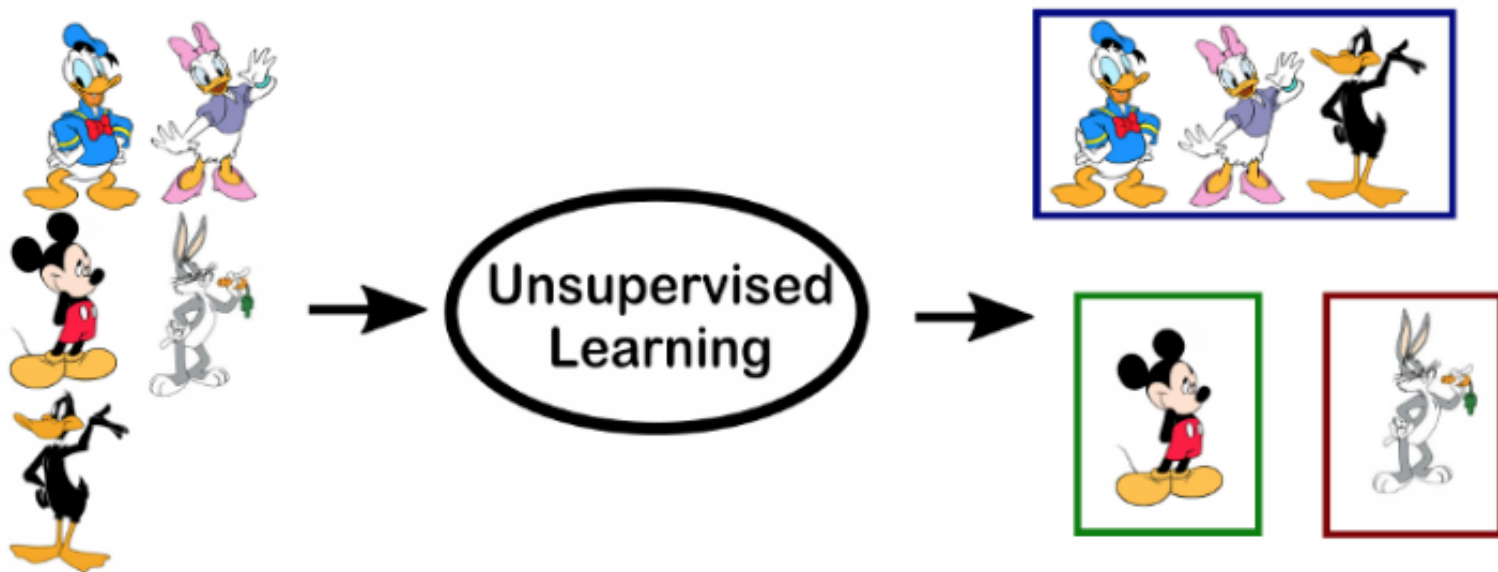
Two Main Types of Machine Learning

- Supervised learning: learn by examples



Two Main Types of Machine Learning

- Unsupervised learning: find structure w/o examples

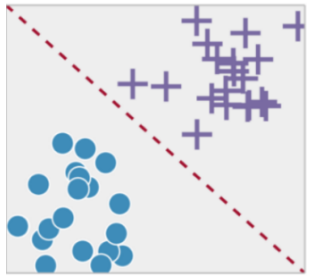
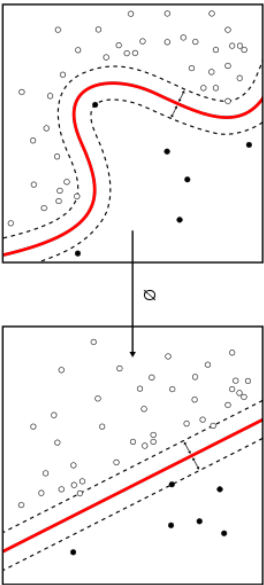
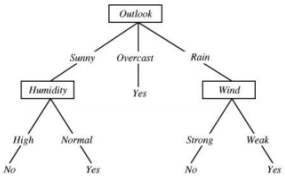
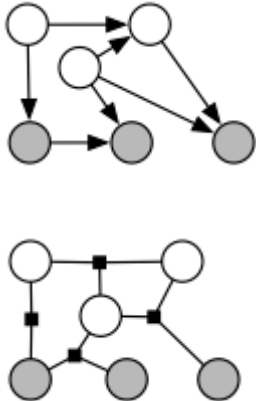
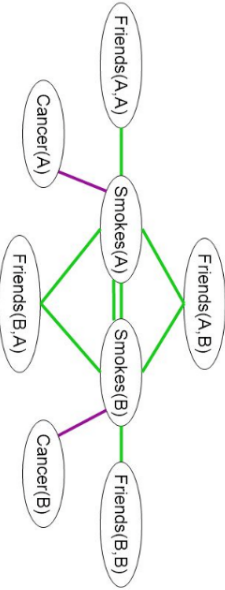
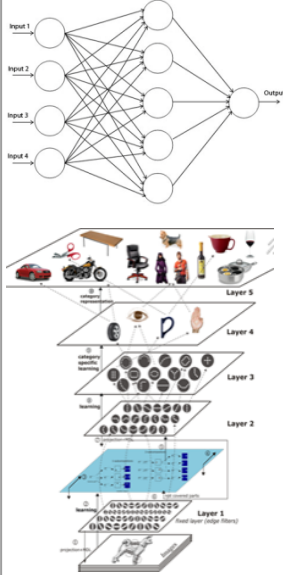


Two Main Types of Machine Learning

- Supervised learning: learn by examples
- Unsupervised learning: find structure w/o examples

	<i>Supervised Learning</i>	<i>Unsupervised Learning</i>
<i>Discrete</i>	classification or categorization	clustering
<i>Continuous</i>	regression	dimensionality reduction

Techniques for Supervised ML

Hyperplanes	Kernel	Tree-based	Graphical Mdl	Logic Prog	Neural Netw
Linear/Logistic regression	SVM	Decision tree, Random forest	Bayes net, CRF	Pr soft logic, Markov logic net	ANN, RNN, CNN
					

Key Lessons for ML [Domingos, 2012]

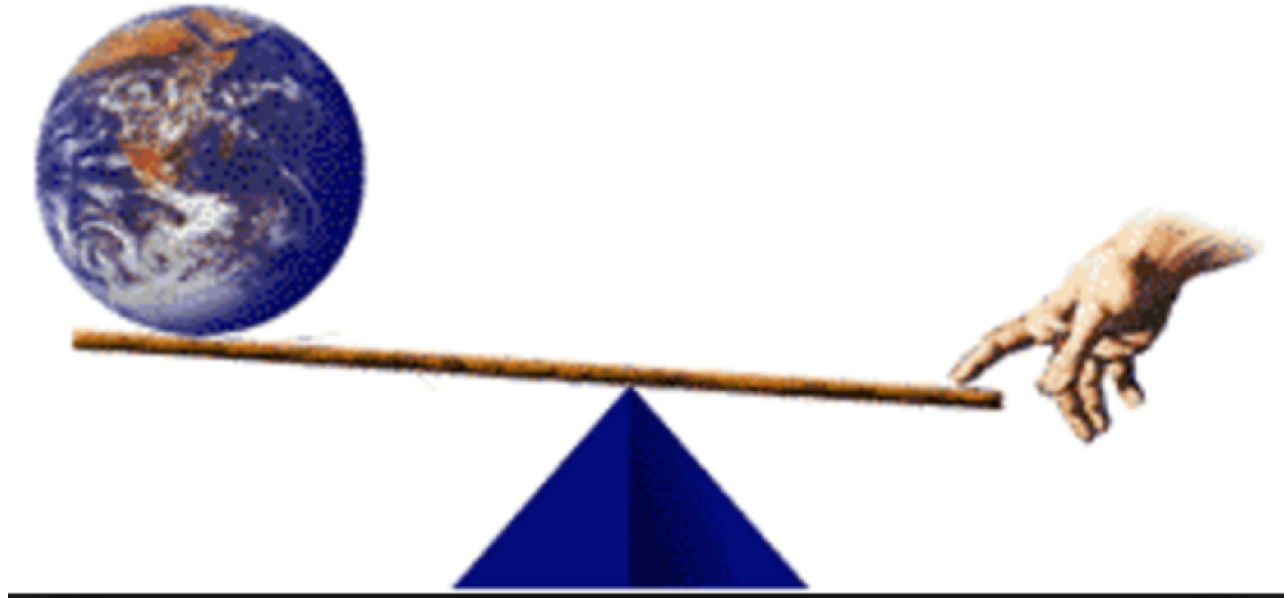
- Learning = Representation + Evaluation + Optimization
- **It's generalization that counts: generalize beyond training examples**
- Data alone is not enough: “no free lunch” theorem--No learner can beat random guessing over all possible functions to be learned
- Intuition fails in high dimensions: “curse of dimensionality”
- **More data beats a cleverer algorithm:** Google showed that after providing 300M images for DL image recognition, no flattening of the learning curve was observed.

DI & ML as Synergy

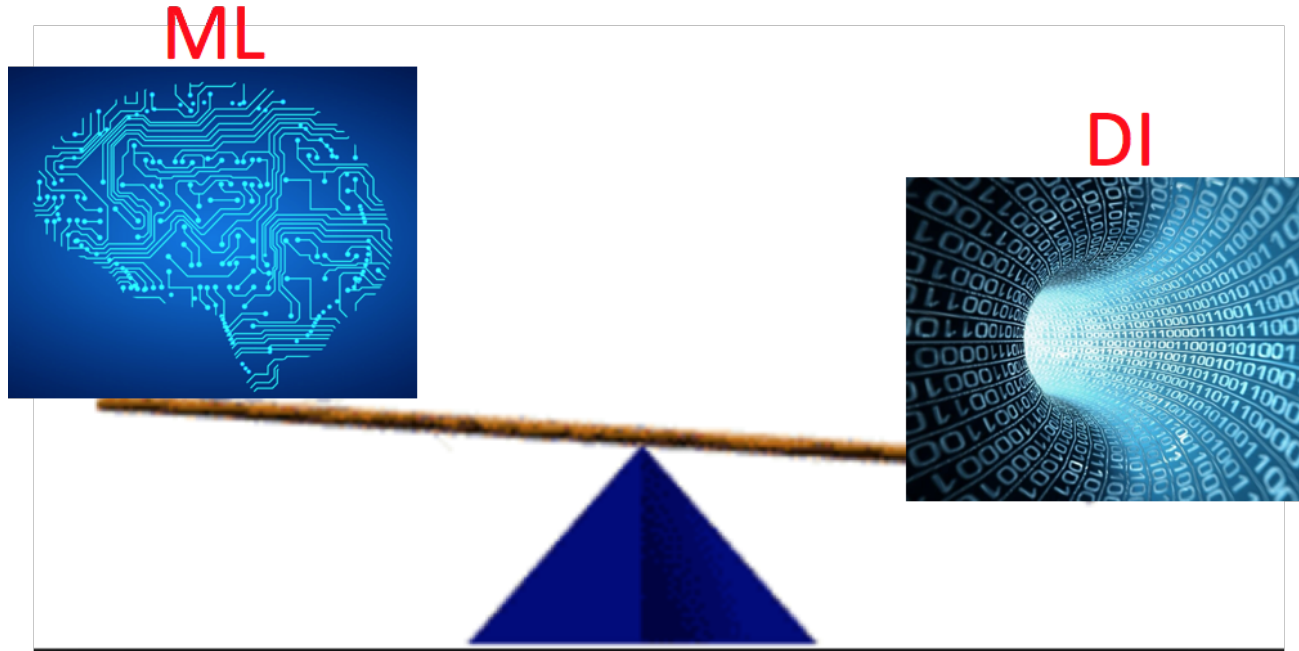
- **ML for effective DI: AUTOMATION, AUTOMATION, AUTOMATION**
 - Automating DI tasks with training data
 - Better understanding of semantics by neural network
- **DI for effective ML: DATA, DATA, DATA**
 - Create large-scale training datasets from different sources
 - Cleaning of data used for training

Give me a Fulscrum, I will Move the Earth

-- Archimedes



Give me a DI funnel, I will Move ML



Many Systems Where DI & ML Leverage Each Other



NELL



MacroBase

QCRI
معهد قطر لبحوث الحوسبة
Qatar Computing Research Institute

جامعة حمد بن خليفة
HAMAD BIN KHALIFA UNIVERSITY



Magellan

HoloClean



snorkel

Dedupe.io



KNOWLEDGE
VAULT



BigGorilla



amperity



product
graph

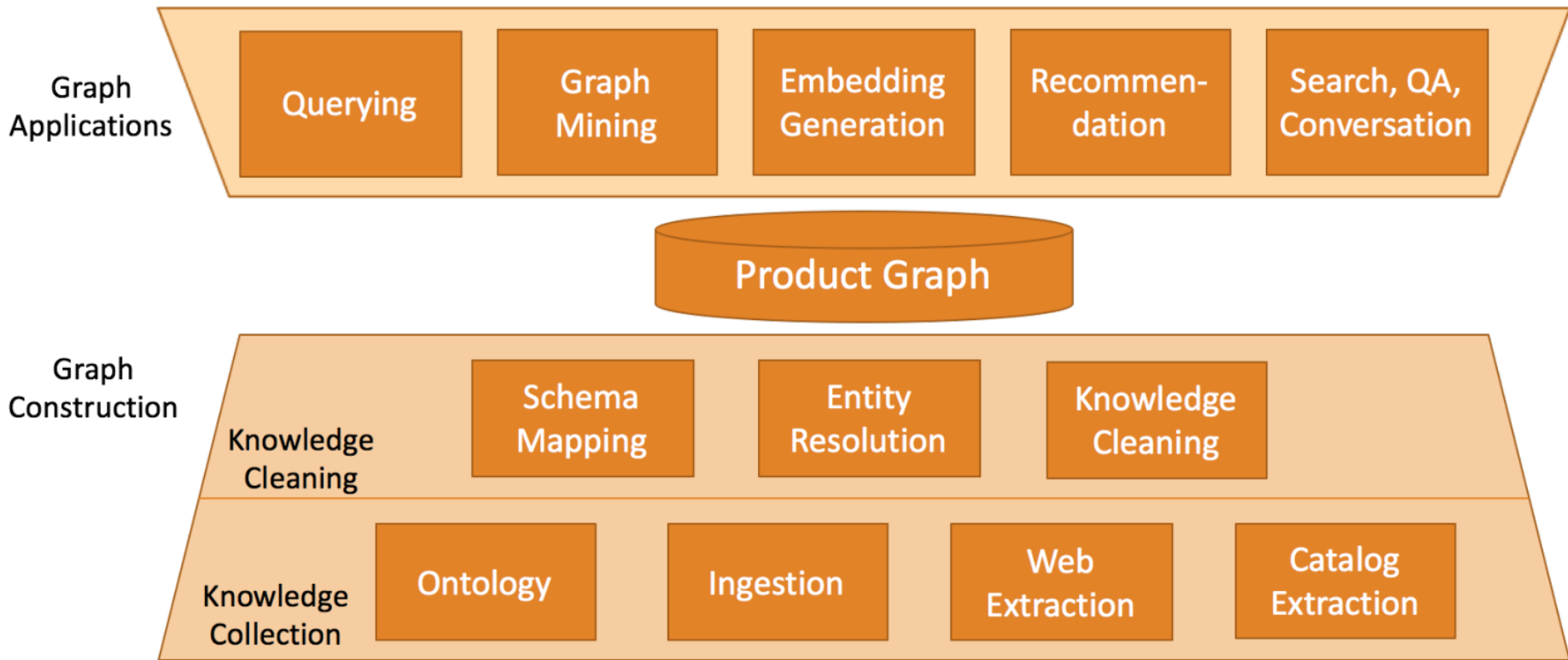
tamr



TRIFACTA

Increasing number of systems both in industry and academia.

Example System: Product Graph [Dong, KDD'18]



Goal of This Tutorial

- **NO-GOALS**

- Present a comprehensive literature review for all topics we are covering

- **GOALS**

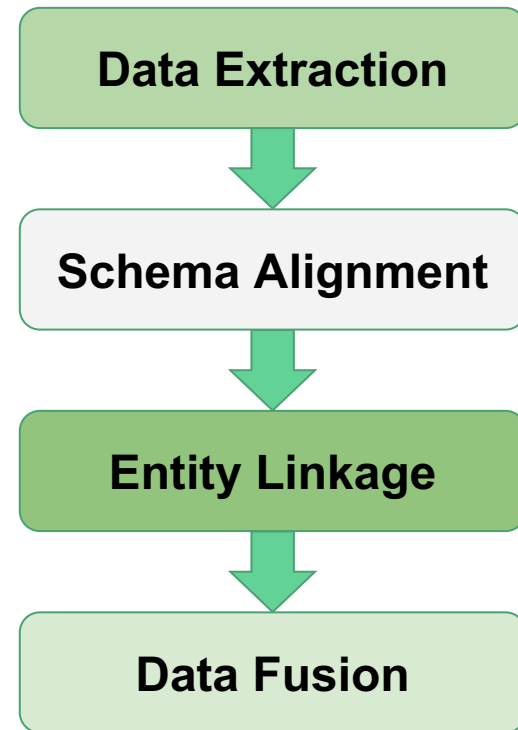
- Present state-of-the-art for DI & ML synergy
- Show how ML has been transforming DI and vice versa
- Give some taste on which tool is working best for which tasks
- Discuss what remains challenging

Outline

- Part I. Introduction
- Part II. ML for DI
- Part III. DI for ML
- Part IV. Conclusions and research directions

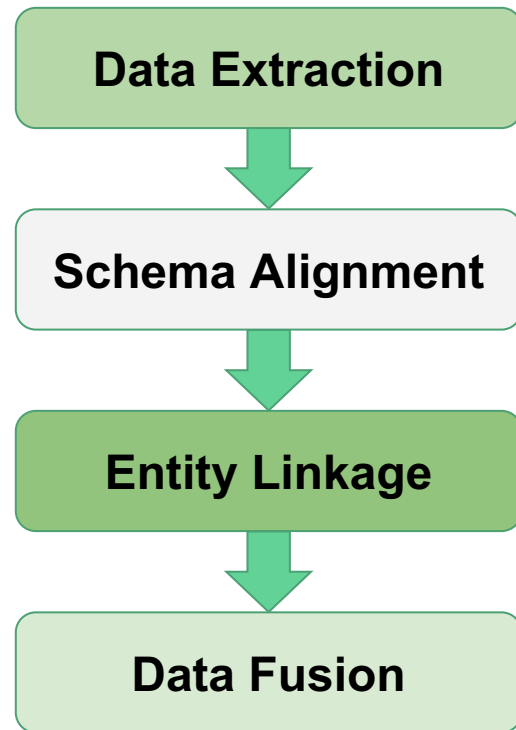
Data Integration Overview

- Entity linkage: linking records to entities; indispensable when different sources exist
- Data extraction: extracting structured data; important when non-relational data exist
- Data fusion: resolving conflicts; necessary in presence of erroneous data
- Schema alignment: aligning types and attributes; helpful when different relational schemas exist



Recipe

- Problem definition
- Brief history
- State-of-the-art ML solutions
- Summary w. a short answer



Theme I. Which ML Model Works Best?



Which ML Model Works Best?

ID	NAME	CLASS	MARK	SEX
1	John Deo	Four	75	female
2	Max Ruin	Three	85	male
3	Arnold	Three	55	male
4	Krish Star	Four	60	female
5	John Mike	Four	60	female
6	Alex John	Four	55	male
7	My John Rob	Fifth	78	male
8	Asruid	Five	85	male
9	Tes Qry	Six	78	male
10	Big John	Four	55	female

Tree-based models

Web tables & Lists

Name and (party) ¹	Term	State of birth	Born
1. Washington (F) ¹	1788		
2. J. Adams (F)	1797		
3. Jefferson (DR)	1801		
4. Madison (DR)	1809		

Free texts

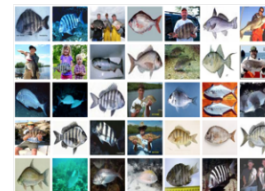
Synopsis [Print](#) [Cite This](#)

Born on April 15, 1452 in Vinci, Italy, Leonardo da Vinci was concerned with the laws of science and nature, which greatly informed his work as a painter, sculptor, inventor and draftsman. His ideas and body of work -- which includes *Virgin of the Rocks*, *The Last Supper*, *Leda and the Swan* and *Mona Lisa* -- have influenced countless artists and made da Vinci a leading light of the Italian Renaissance.

??

SCENE FROM "DAN'L DRUCE."

This interesting domestic drama, by Mr. W. S. Gilbert, has continued to engage the sympathies of a nightly efficient audience at the Haymarket Theatre, where it has now been represented more than sixty times. Its subject and character were described by us, in the ordinary report of theatrical novelties, about two months ago. Our readers will probably not need to be reminded that the hero of the story, Dan'l Druce, the blacksmith, is a solitary recluse dwelling on the coast of Norfolk, where his lone cottage is visited by fugitives from party vengeance during the civil wars of the Commonwealth. His hoard of money is stolen; but a different sort of treasure, a helpless female infant; is left by some mysterious agency, and may be accepted, as in George Eliot's tale of "Silas Marner," for a living gift to the sad-hearted misanthrope, far better than riches. In this spirit, at least, he is content to receive the precious human charge; and so to those who would remove it from his home, Dan'l Druce here makes answer with the solemn exclamation, "Touch not the Lord's gift!" This character is well acted by Mr. Hermann Vezin.



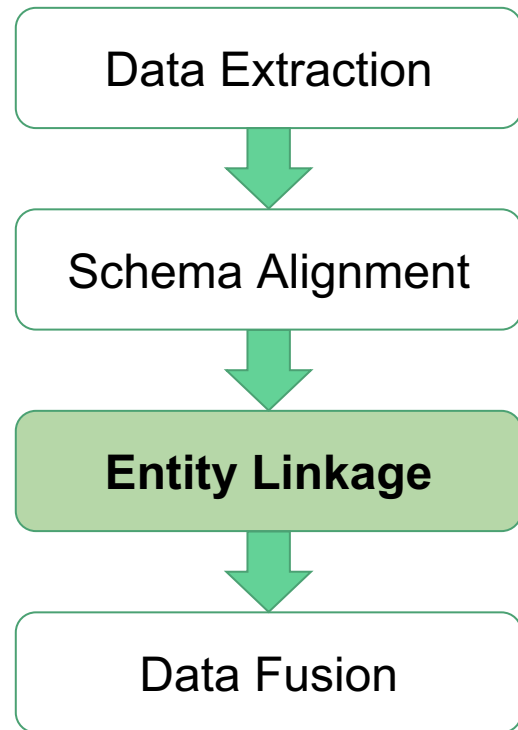
Neural network

Theme II. Does Supervised Learning Apply to DI?

- Supervised learning has made a big splash recently in many fields
- However, it is hard to bluntly apply supervised learning to DI tasks
 - Our goal is to integrate data from many different data sources in different domains
 - The different sources present different data features and distributions
 - Collecting training labels for each source is a huge cost

Outline

- Part I. Introduction
- Part II. ML for DI
 - ML for entity linkage
 - ML for data extraction
 - ML for data fusion
 - ML for schema alignment
- Part III. DI for ML
- Part IV. Conclusions and research direction



What is Entity Linkage?

- Definition: Partition a given set R of records, such that each partition corresponds to a distinct real-world entity.

Are they the same entity?

IMDB



Anahí

Actress | Music Department | Soundtrack

Anahi was born in Mexico. She's had roles in *Tuy Yo*, in which she played a 17 year old girl while she was 13, and *Vivo Por Elena*, in which she played Talita, a naive and innocent teenager. Anahi lives with her mother and sister name Marychelo. She hopes to become a fashion designer one day, and is currently pursuing a career in singing.

[See full bio »](#)

Born: May 14, 1982 in Mexico City, Distrito Federal, Mexico

[More at IMDbPro](#) »

📞 Contact Info: [View manager](#)

 **SEE RANK**

WikiData

Anahí Puente (Q169461)

Mexican singer-songwriter and actress

Mia

▼ In more languages [Configure](#)

Language	Label	Description
English	Anahí Puente	Mexican singer-songwriter and actress
Chinese	阿纳希·普恩特	No description defined
Spanish	Anahí Puente	Cantante, compositora y actriz mexicana

date of birth

7 November 1983

▼ 1 reference

imported from

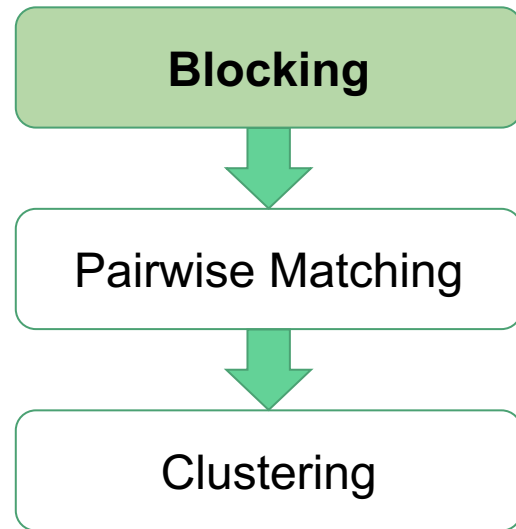
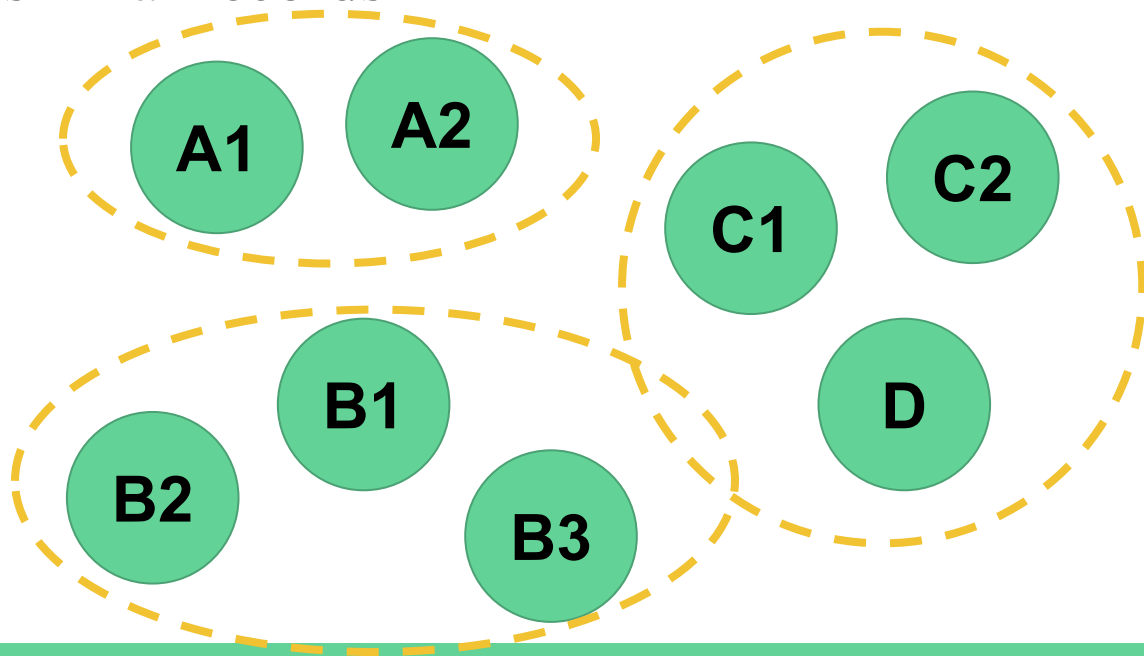
Italian Wikipedia

+ [add reference](#)

+ add value

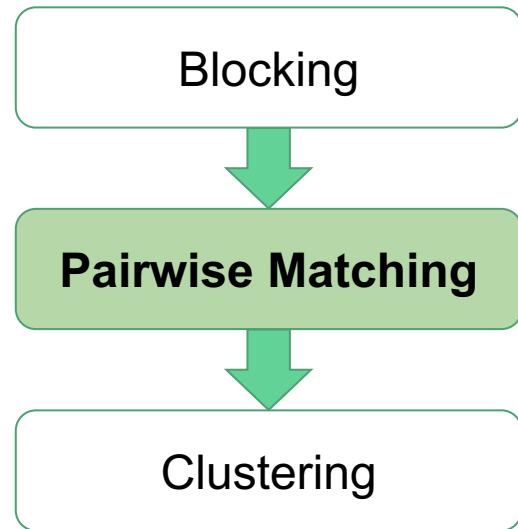
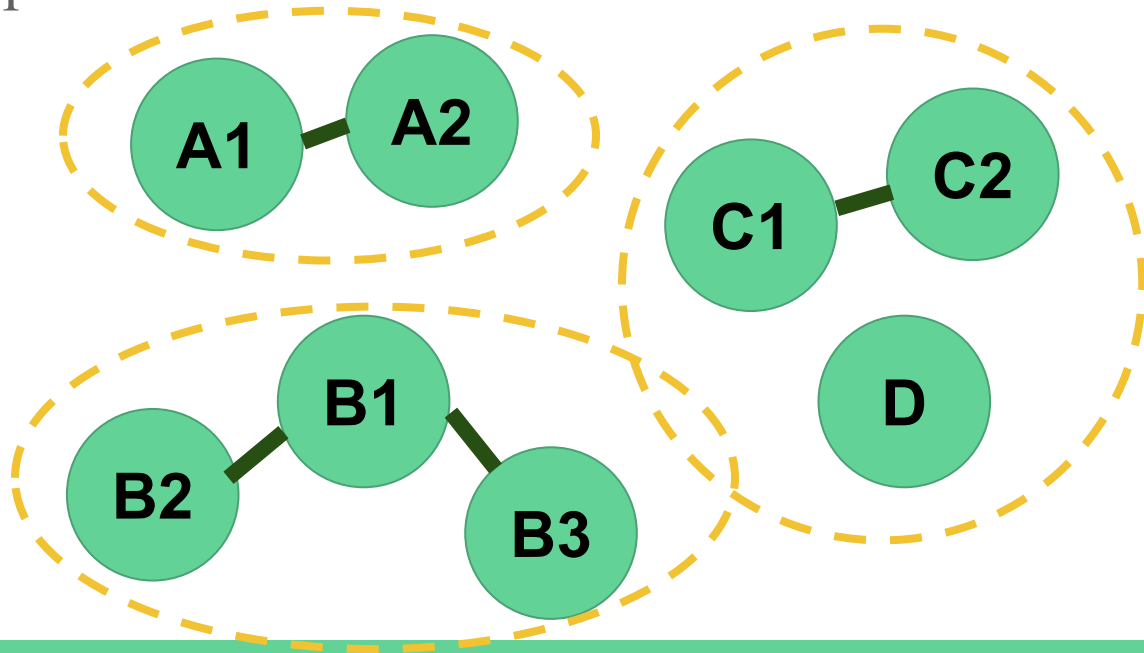
Three Steps in Entity Linkage

- **Blocking:** efficiently create small blocks of similar records



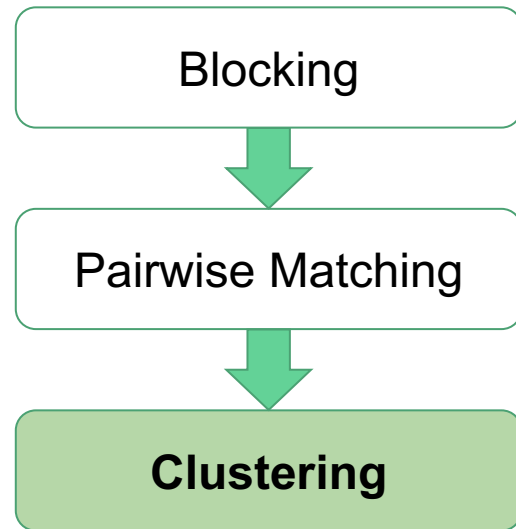
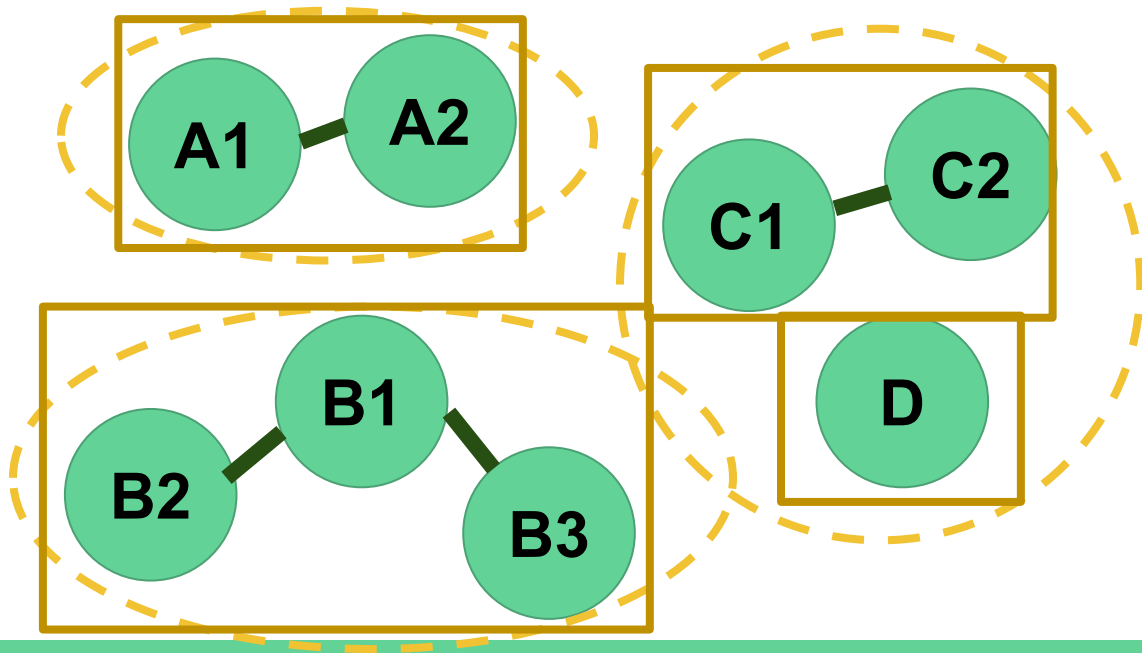
Three Steps in Entity Linkage

- **Pairwise matching:** compare all record pairs in a block

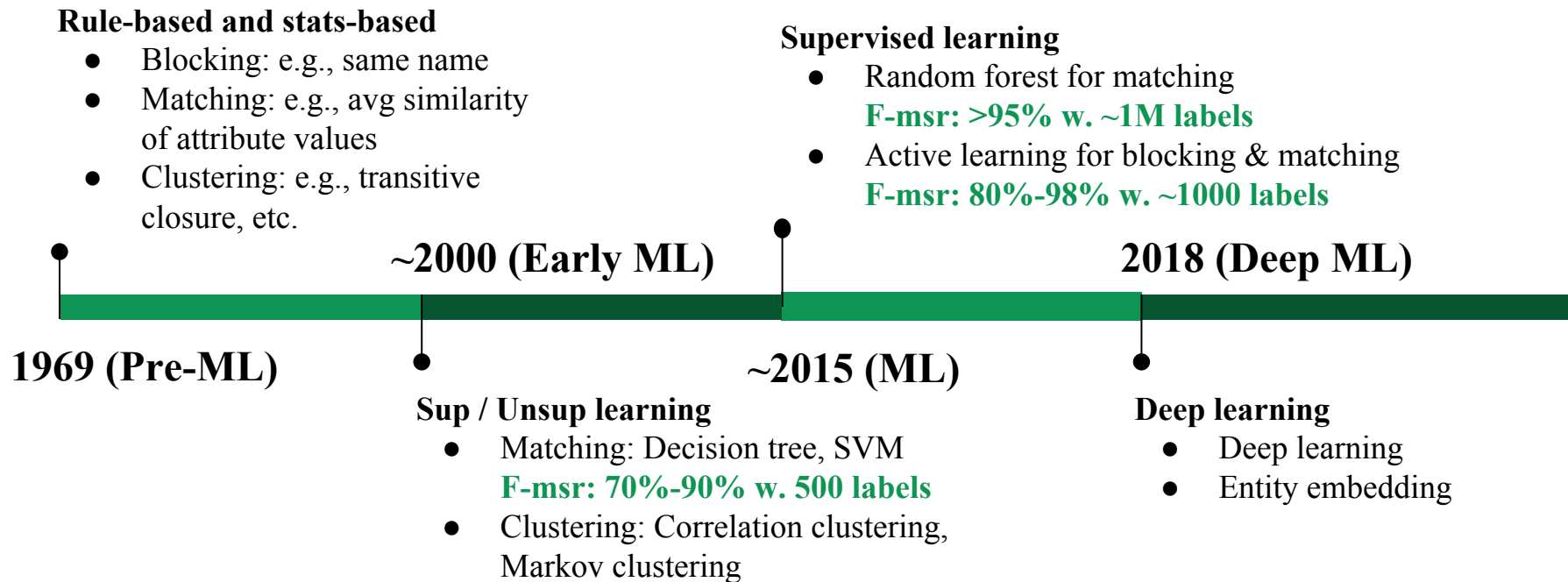


Three Steps in Entity Linkage

- **Clustering:** group records into entities



50 Years of Entity Linkage



Rule-Based Solution

Rule-based and stats-based

- Blocking: e.g., same name
- Matching: e.g., avg similarity of attribute values
- Clustering: e.g., transitive closure, etc.

● [Fellegi and Sunter, 1969]

- Match: $\text{sim}(r, r') > \theta_h$
- Unmatch: $\text{sim}(r, r') < \theta_l$
- Possible match:
 $\theta_l < \text{sim}(r, r') < \theta_h$

1969 (Pre-ML)

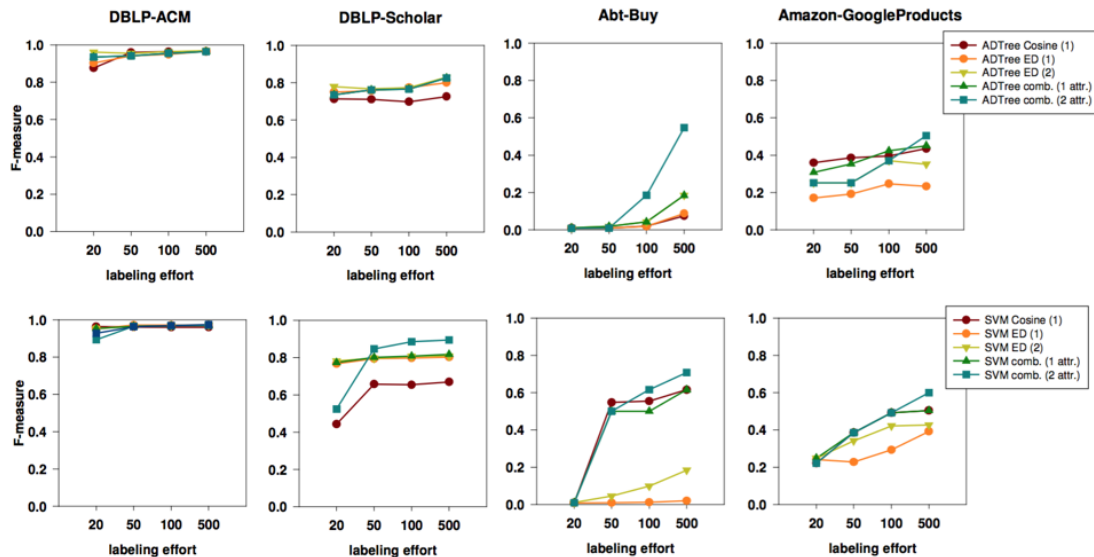
Early ML Models

- [Köpcke et al, VLDB'10]

~2000 (Early ML)

Sup / Unsup learning

- Matching: Decision tree, SVM
F-msr: 70%-90% w. 500 labels
- Clustering: Correlation clustering, Markov clustering



State-of-the-Art ML Models [Dong, KDD'18]

Supervised learning

- Random forest for matching
F-msr: >95% w. ~1M labels
- AL for blocking & matching
F-msr: 80%-98% w. ~1000 labels

~2015 (ML)

- Features: attribute similarity measured in various ways. E.g.,
 - string sim: Jaccard, Levenshtein
 - number sim: absolute diff, relative diff
- ML models on Freebase vs. IMDb
 - Logistic regression: Prec=0.99, Rec=0.6
 - Random forest: Prec=0.99, Rec=0.99

State-of-the-Art ML Models [Dong, KDD'18]

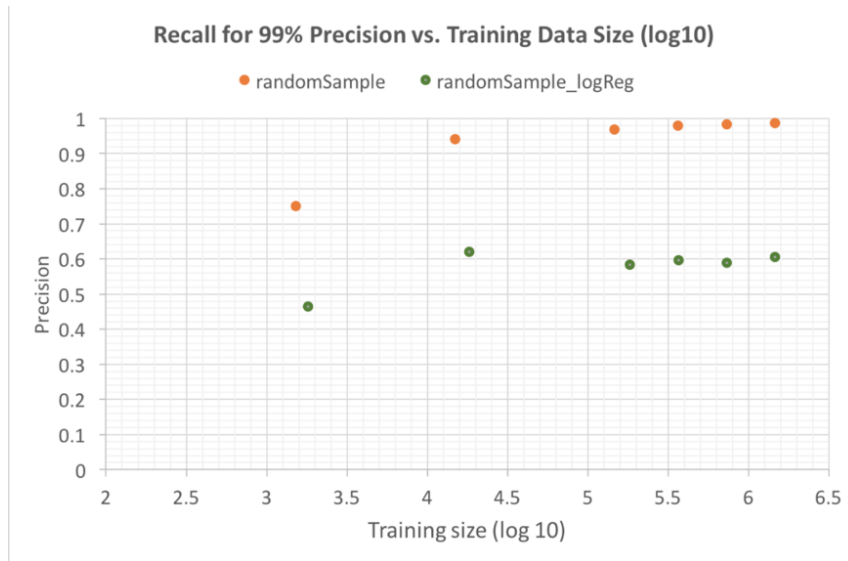
Supervised learning

- Random forest for matching
F-msr: >95% w. ~1M labels
- AL for blocking & matching
F-msr: 80%-98% w. ~1000 labels

~2015 (ML)

● Expt 1. IMDb vs. Freebase

- Logistic regression: Prec=0.99, Rec=0.6
- Random forest: Prec=0.99, Rec=0.99



State-of-the-Art ML Models [Dong, KDD'18]

Supervised learning

- Random forest for matching
F-msr: >95% w. ~1M labels
- AL for blocking & matching
F-msr: 80%-98% w. ~1000 labels

~2015 (ML)

- Features: attribute similarity measured in various ways. E.g.,
 - name sim: Jaccard, Levenshtein
 - age sim: absolute diff, relative diff
- ML models on Freebase vs. IMDb
 - Logistic regression: Prec=0.99, Rec=0.6
 - Random forest: Prec=0.99, Rec=0.99
 - XGBoost: marginally better, but sensitive to hyper-parameters

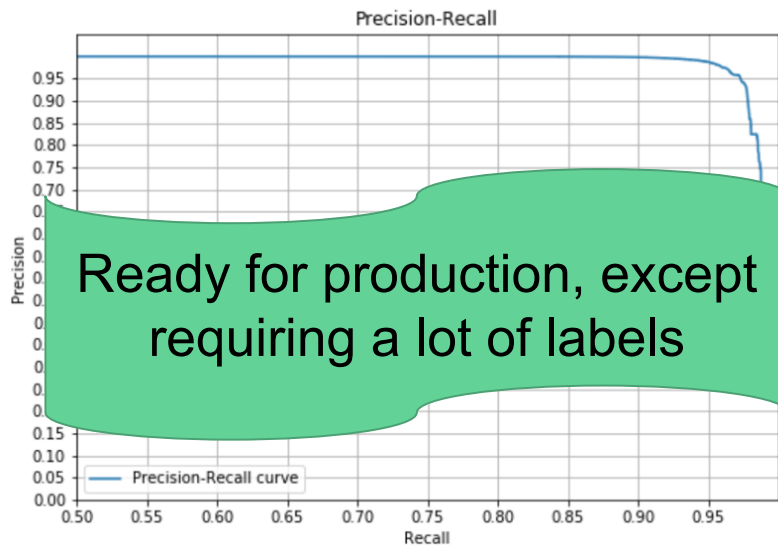
State-of-the-Art ML Models [Dong, KDD'18]

Supervised learning

- Random forest for matching
F-msr: >95% w. ~1M labels
- AL for blocking & matching
F-msr: 80%-98% w. ~1000 labels

~2015 (ML)

- Expt 2. IMDb vs. Amazon movies
 - 200K labels, ~150 features
 - Random forest: Prec=0.98, Rec=0.95





Magellan

State-of-the-Art ML Models [Das et al., SIGMOD'17]

Supervised learning

- Random forest for matching
F-msr: >95% w. ~1M labels
- AL for blocking & matching
F-msr: 80%-98% w. ~1000 labels

~2015 (ML)

- Falcon: apply active learning both for blocking and for matching; ~1000 labels

Dataset	Accuracy (%)			Cost (# Questions)
	P	R	F_1	
Products	90.9	74.5	81.9	\$57.6 (960)
Songs	96.0	99.3	97.6	\$54.0 (900)
Citations	92.0	98.5	95.2	\$65.5 (1087)

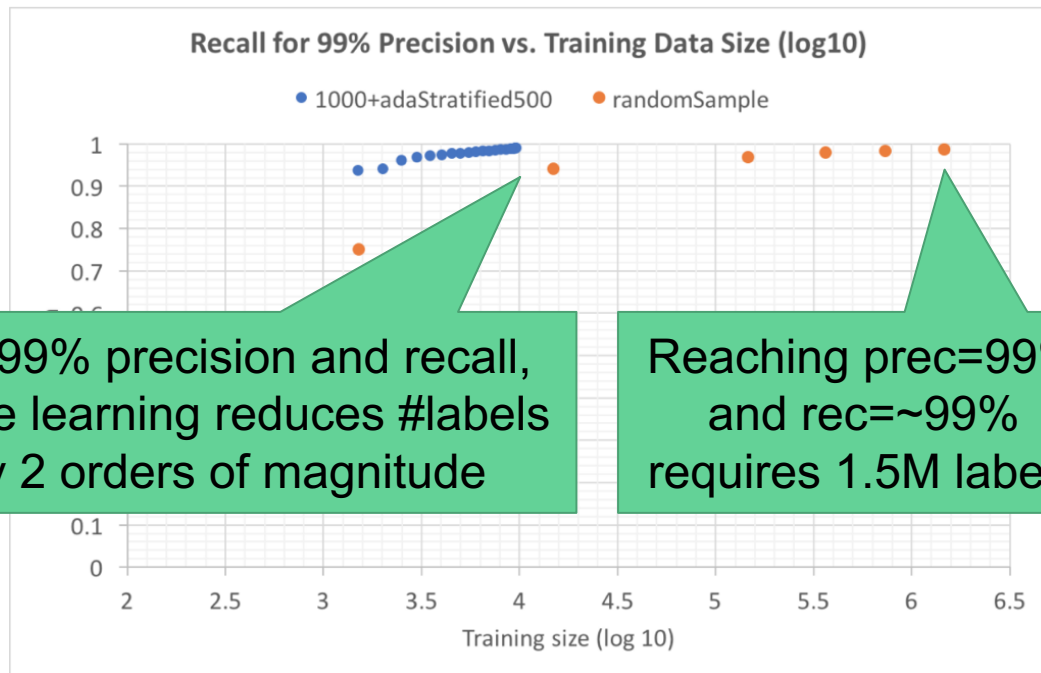
State-of-the-Art ML Models [Dong, KDD'18]

Supervised learning

- Random forest for matching
F-msr: >95% w. ~1M labels
- AL for blocking & matching
F-msr: 80%-98% w. ~1000 labels

~2015 (ML)

- Apply active learning to minimize #labels



Deep Learning Models [Mudgal et al., SIGMOD'18]



Magellan

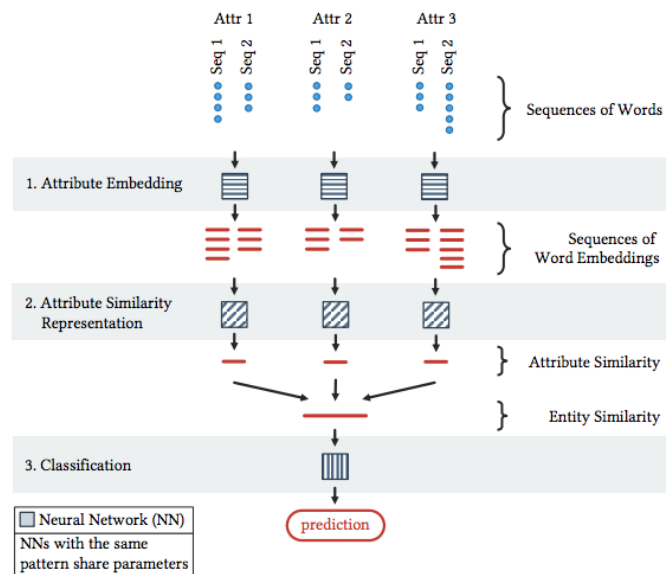
Check-out at poster session
on Wednesday!
Code at: [deepmatcher.ml](https://github.com/mudgal/deepmatcher.ml)

- Bi-RNN w. attention
- Similar performance for structured data;
Significant improvement on texts and dirty data

2018 (Deep ML)

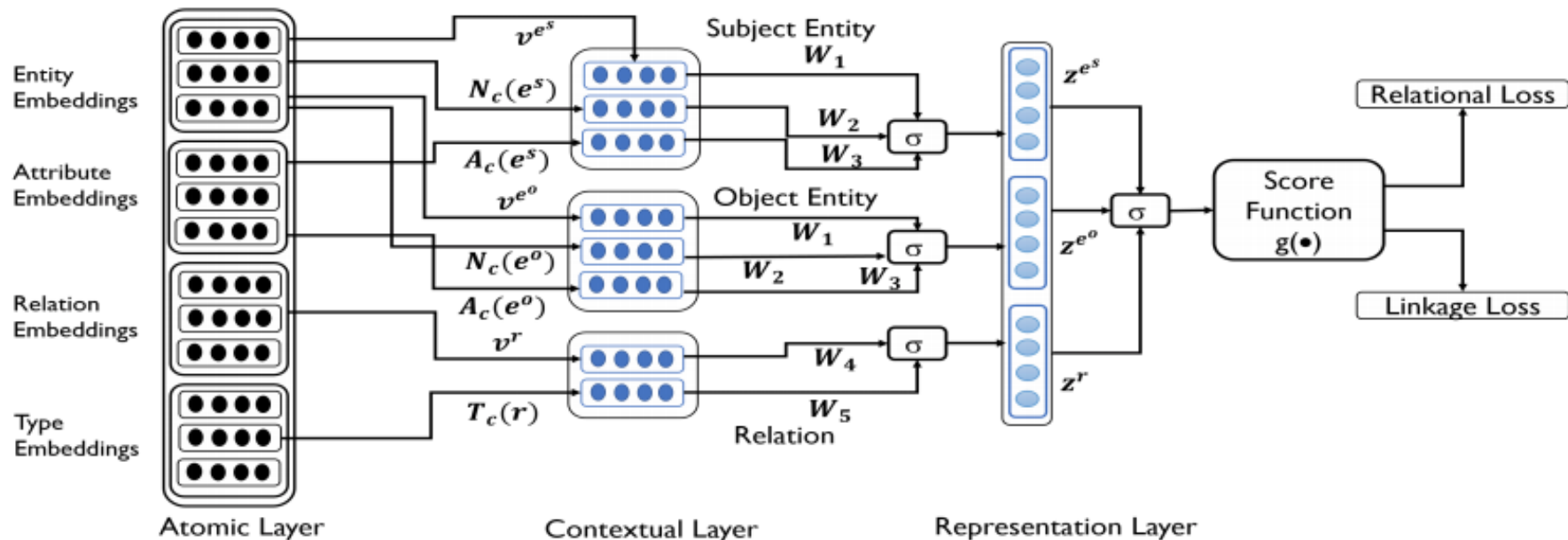
Deep learning

- Deep learning
- Entity embedding



Deep Learning Models [Trivedi et al., ACL'18]

- LinkNBed: Generate embeddings for entities as in knowledge embedding



Deep Learning Models [Trivedi et al., ACL'18]

- LinkNBed: Generate embeddings for entities as in knowledge embedding
- Performance better than previous knowledge embedding methods, but not comparable to random forest
- Enable linking different types of entities

2018 (Deep ML)

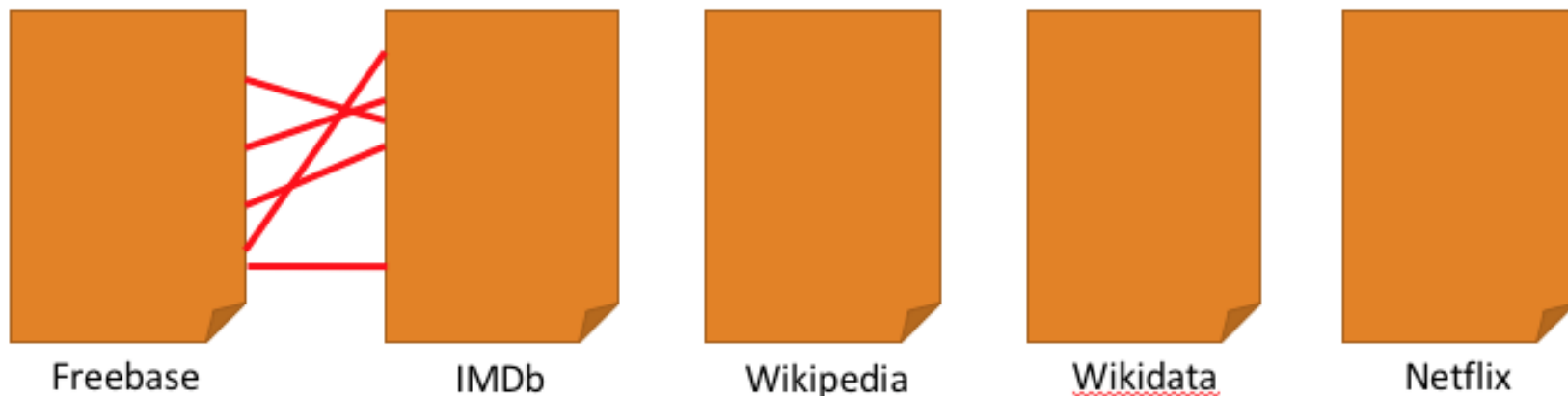


Deep learning

- Deep learning
- Entity embedding

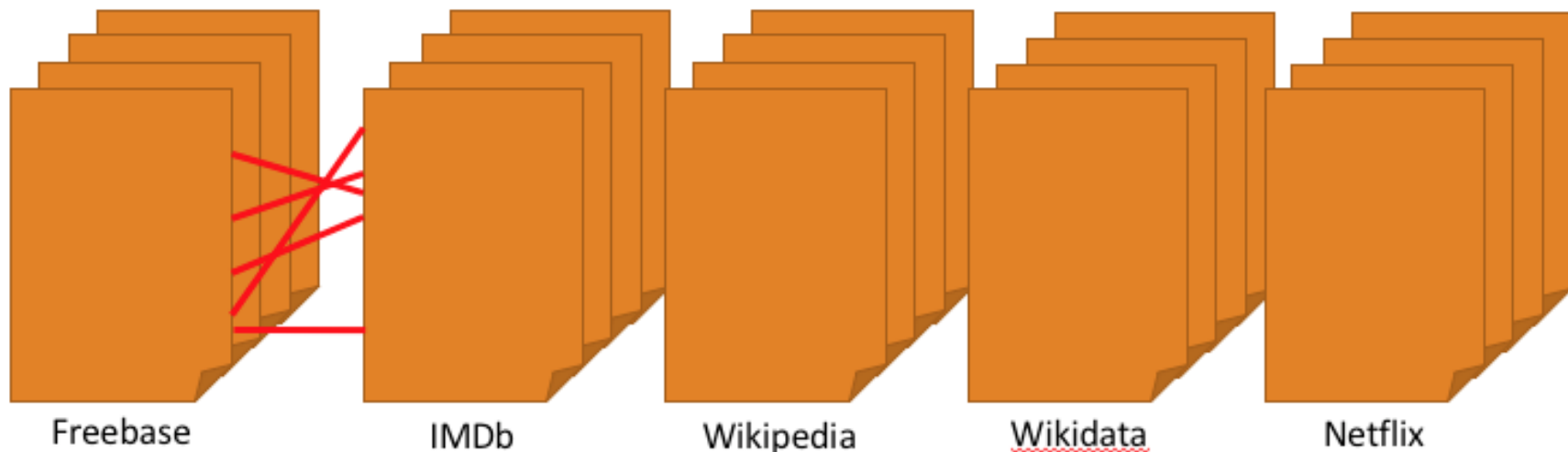
Challenges in Applying ML on EL

- How can we obtain abundant training data for many types, many sources, and dynamically evolving data??
- From two sources to multiple sources



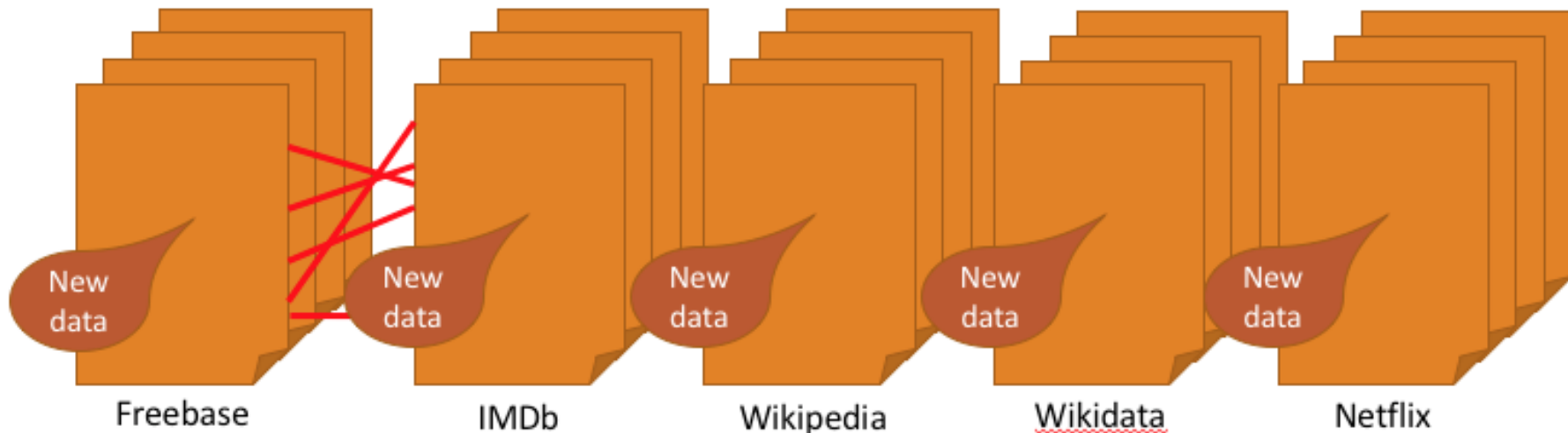
Challenges in Applying ML on EL

- How can we obtain abundant training data for many types, many sources, and dynamically evolving data??
- From one entity type to multiple types



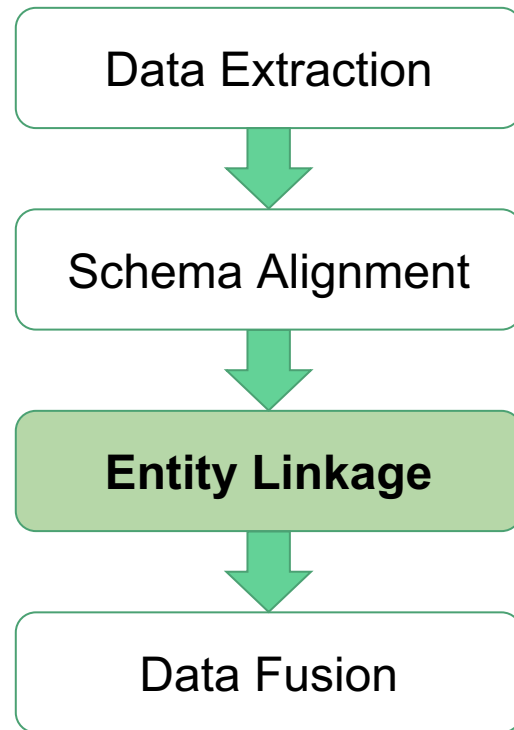
Challenges in Applying ML on EL

- How can we obtain abundant training data for many types, many sources, and dynamically evolving data??
- From static data to dynamic data



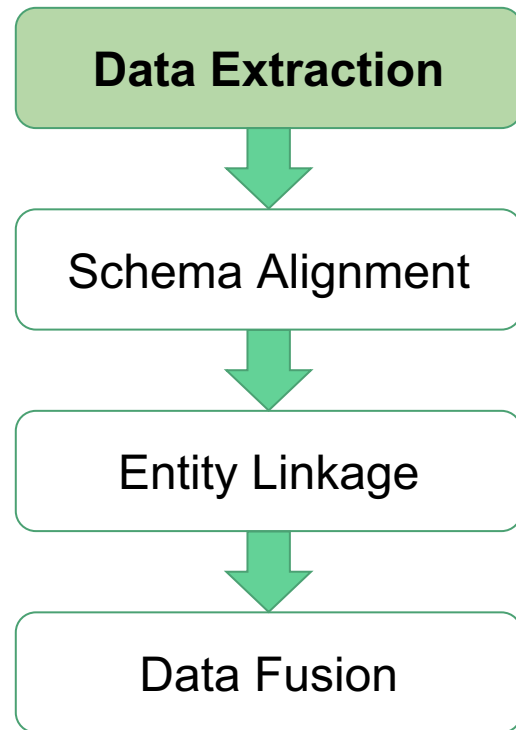
Recipe for Entity Linkage

- Problem definition: **Link references to the same entity**
- Short answers
 - **RF w. attribute-similarity features**
 - **DL to handle texts and noises**
 - **End-to-end solution is future work**



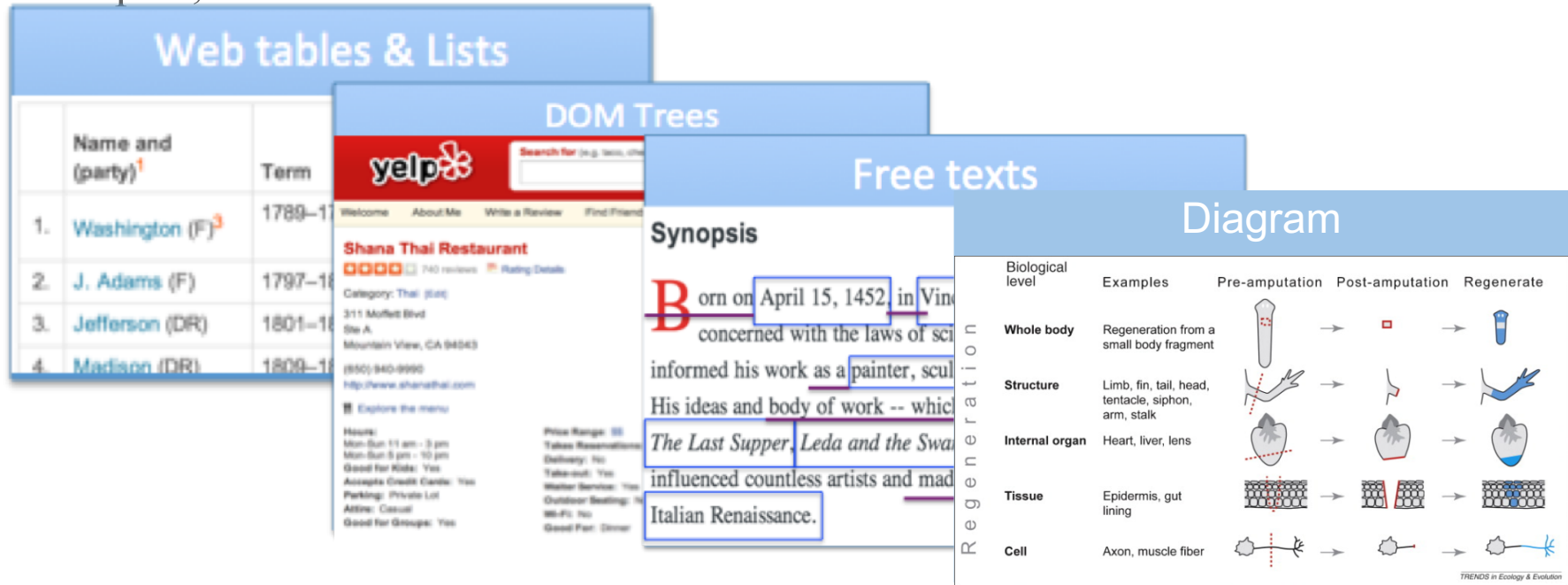
Outline

- Part I. Introduction
- Part II. ML for DI
 - ML for entity linkage
 - ML for data extraction
 - ML for data fusion
 - ML for schema alignment
- Part III. DI for ML
- Part IV. Conclusions and research direction



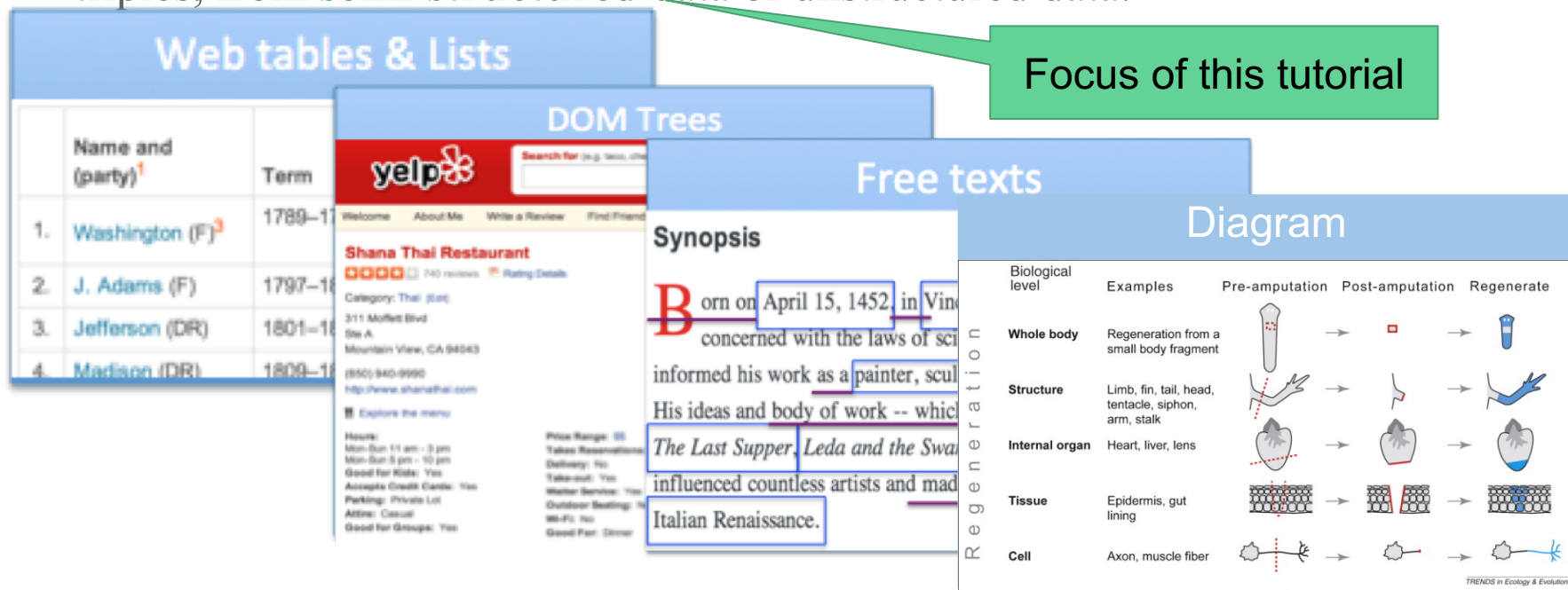
What is Data Extraction?

- Definition: Extract structured information, e.g., (entity, attribute, value) triples, from semi-structured data or unstructured data.



What is Data Extraction?

- Definition: Extract structured information, e.g., (entity, attribute, value) triples, from **semi-structured** data or unstructured data.



Three Types of Data Extraction

- **Closed-world extraction:** align to existing entities and attributes; e.g., (ID_Obama, place_of_birth, ID_USA)
- **ClosedIE:** align to existing attributes, but extract new entities; e.g., (“Xin Luna Dong”, place_of_birth, “China”)
- **OpenIE:** not limited by existing entities or attributes; e.g., (“Xin Luna Dong”, “was born in”, “China”), (“Luna”, “is originally from”, “China”)

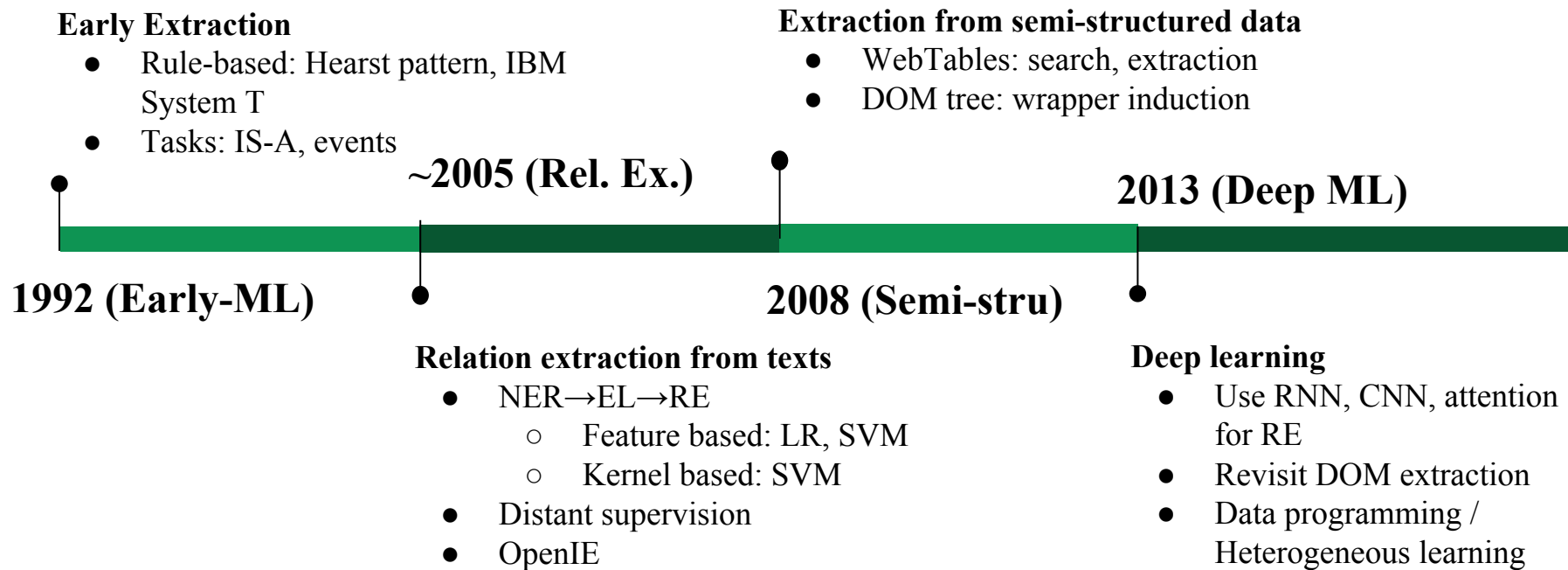
Three Types of Data Extraction

- **Closed-world extraction:** align to existing entities and attributes; e.g., (ID_Obama, place_of_birth, ID_USA)
- **ClosedIE:** align to existing attributes, but extract new entities; e.g., (“Xin Luna Dong”, place_of_birth, “China”)
- **OpenIE:** not limited by existing entities or attributes; e.g., (“Xin Luna Dong”, “was born in”, “China”), (“Luna”, “is originally from”, “China”)

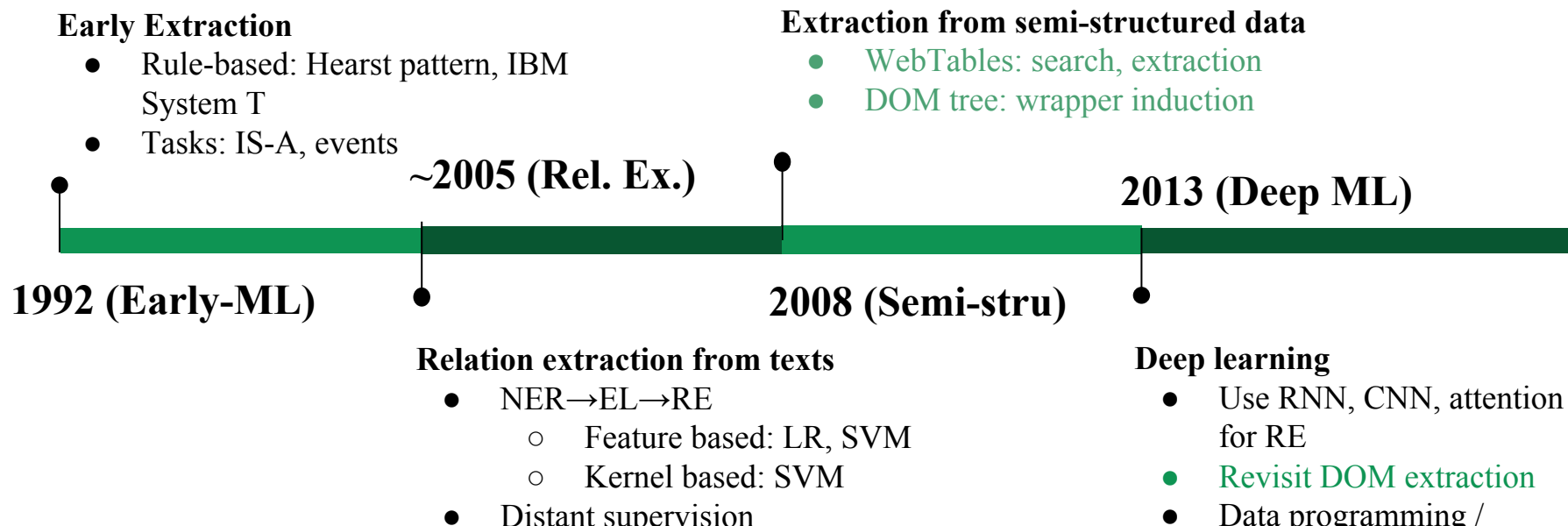


Focus of this tutorial

35 Years of Data Extraction



35 Years of Data Extraction

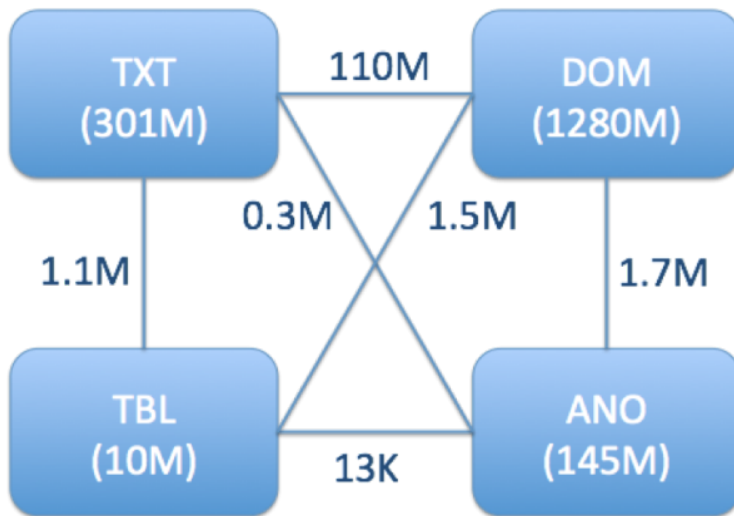


Come to our VLDB tutorial for text extraction and OpenIE!!

Why Semi-Structured Data?

- Knowledge Vault @ Google showed big potential from DOM-tree extraction
[Dong et al., KDD'14][Dong et al., VLDB'14]

Accu	Accu (conf $\geq .7$)
0.36	0.52



Accu	Accu (conf $\geq .7$)
0.43	0.63
0.09	0.62

Wrapper Induction--Vertex [Gulhane et al., ICDE'11]

Title **Genre** **Release Date**

Runtime

PG | 1h 50min | Action, Drama, Romance | 16 May 1986 (USA)

Watch Now
From \$2.99 (SD) on Amazon Video

As students at the United States Navy's elite fighter weapons school compete to be best in the class, one daring young pilot learns a few things from a civilian instructor that are not taught in the classroom.

Director Tony Scott
Writers Jim Cash, Jack Epps Jr. | 1 more credit »
Stars Tom Cruise, Tim Robbins, Kelly McGillis | See full cast & crew »

Actors

Metascore 50 From metacritic.com | Reviews 404 user | 173 critic | Popularity 404 (71)

Extracted relationships

- (Top Gun, type.object.name, "Top Gun")
- (Top Gun, film.film.genre, Action)
- (Top Gun, film.film.directed_by, Tony Scott)
- (Top Gun, film.film.starring, Tom Cruise)
- (Top Gun, film.film.runtime, "1h 50min")
- (Top Gun, film.film.release_Date_s, "16 May 1986")

Wrapper Induction--Vertex [Gulhane et al., ICDE'11]

- Solution: find XPaths from DOM Trees

Filmography		Show all	Show by...	Edit
Jump to: Actor Producer Soundtrack Director Writer Thanks Self Archive footage				
Actor (46 credits)		Hide ▲		
Top Gun: Maverick (pre-production)		2019		
Maverick				
M:I - Mission Impossible (filming)		2018		
Ethan Hunt				
American Made (completed)		2017		
Barry Seal				
Luna Park (announced)				
The Mummy		2017		
Nick Morton				
Jack Reacher: Never Go Back		2016		
Jack Reacher				
Mission: Impossible - Rogue Nation		2015		
Ethan Hunt				
Edge of Tomorrow		2014		
Cage				
Oblivion		2013/1		
Jack				
Jack Reacher		2012		
Reacher				
Rock of Ages		2012		
Stacey Jaxx				
Mission: Impossible - Ghost Protocol		2011		
Ethan Hunt				
Knight and Day		2010		
Roy Miller				
Valkyrie		2008		
Colonel Claus von Stauffenberg				
Tropic Thunder		2008		

```
<div id="filmography"> == $0
  <div id="filmo-head-actor" class="head" data-category="actor" onclick=
    "toggleFilmoCategory(this);"></div>
  <div class="filmo-category-section">
    <div class="filmo-row odd" id="actor-tt1745960">
      <span class="year_column">
        &nbsp;2019
      </span>
      <b>
        <a href="/title/tt1745960/?ref=nm_flg_act_1">Top Gun: Maverick</a>
      </b>
      "
      <a href="/r/legacy-inprod-name/title/tt1745960" class="in_production">pre-
        production</a>
      "
      <br>
      <a href="/character/ch0005702/?ref=nm_flg_act_1">Maverick</a>
    </div>
    <div class="filmo-row even" id="actor-tt4912910"></div>
    <div class="filmo-row odd" id="actor-tt3532216"></div>
    <div class="filmo-row even" id="actor-tt1123441"></div>
    <div class="filmo-row odd" id="actor-tt2345759">
      <span class="year_column">
        &nbsp;2017
      </span>
      <b>
        <a href="/title/tt2345759/?ref=nm_flg_act_5">The Mummy</a>
      </b>
      <br>
      <a href="/character/ch0573416/?ref=nm_flg_act_5">Nick Morton</a>
    </div>
    <div class="filmo-row even" id="actor-tt3393786"></div>
    <div class="filmo-row odd" id="actor-tt2381249"></div>
    <div class="filmo-row even" id="actor-tt1631867"></div>
    <div class="filmo-row odd" id="actor-tt1483013"></div>
    <div class="filmo-row even" id="actor-tt0790724"></div>
    <div class="filmo-row odd" id="actor-tt1336080"></div>
```

Wrapper Induction--Vertex [Gulhane et al., ICDE'11]

- Challenge: slight variations from page to page

```
/html/body/div[1]/div/div[4]/div[3]/div[3]/div[3]/div[3]/div[4]/div[26]/b/a  
/html/body/div[1]/div/div[4]/div[3]/div[3]/div[3]/div[3]/div[2]/div[10]/b/a
```

Figure 2: Example of XPath expressions corresponding to the *acted in* predicate on two IMDb pages. They differ at two node indices, and the second path corresponds to the *producer of* predicate from the first page.

Wrapper Induction--Vertex [Gulhane et al., ICDE'11]

Identify representative webpages for annotation



Learn

Web site
sample pages

Cluster
Pages

Sample pages

Annotate
Pages

Annotations

Learn XSLT Rules

One website may use
multiple templates
(Unsupervised-clustering)

Sample pages

Monitor
Rules

Changed sites

Combine attr features
and textual features to
find a general XPath
(LR)

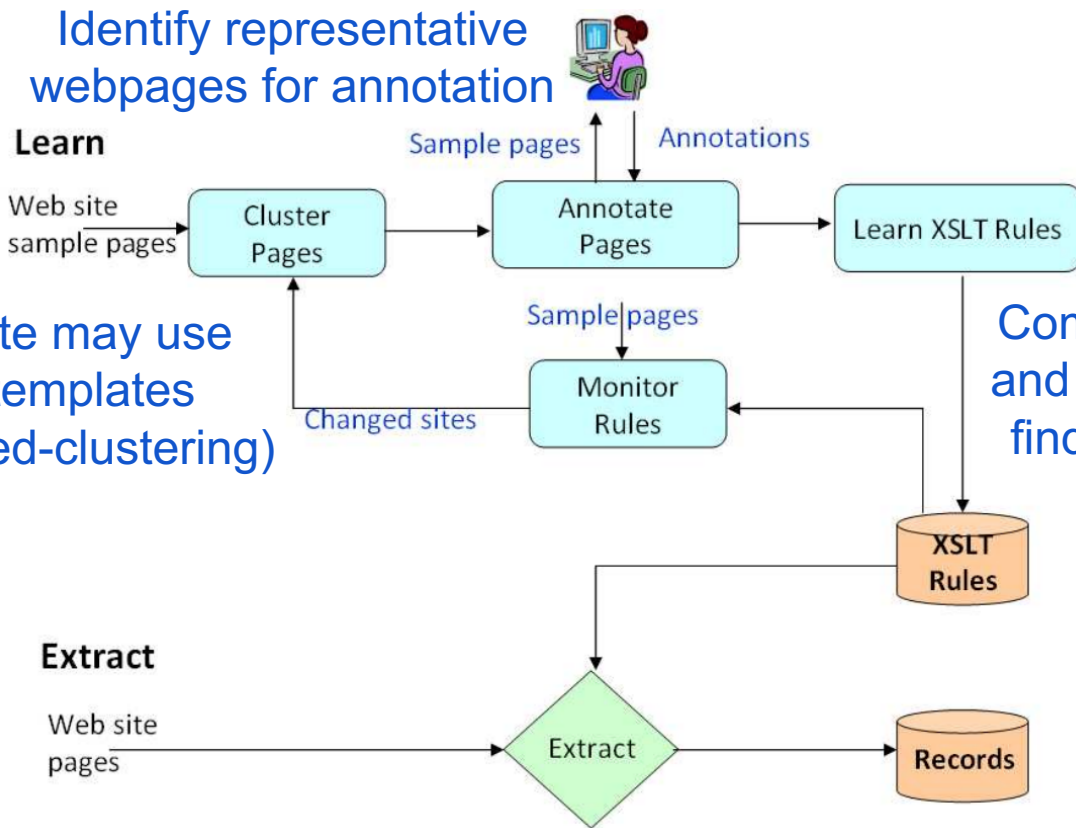
XSLT
Rules

Extract

Web site
pages

Extract

Records



Wrapper Induction--Vertex [Gulhane et al., ICDE'11]

- **Sample learned XPathS on IMDb**

- `//*[@itemprop="name"]`
- `//*[@class="bp_item bp_text_only"]/*/*/*[@class="bp_heading"]`
- `//*[following-sibling::*[position()=3][@class="subheading"]/*[following-sibling::*[position()=1][@class="attribute"]]`
- `//*[preceding-sibling::node()[normalize-space(.)!=""][text()="Language:"]`

Ensure high recall

Ensure high precision

Distantly Supervised Extraction

- **Annotation-based extraction**
 - Pros: high precision and recall
 - Cons: does not scale--annotation per cluster per website
- **Distantly-supervised extraction**
 - Step 1. Use seed data to automatically annotate
 - Step 2. Use the (noisy) annotations for training
 - E.g., DeepDive, Knowledge Vault

Distant Supervision [Mintz et al., ACL'09]

Corpus Text

Bill Gates founded Microsoft in 1975.
Bill Gates, founder of Microsoft, ...
Bill Gates attended Harvard from ...
Google was founded by Larry Page ...

Training Data

Freebase

(Bill Gates, Founder, Microsoft)
(Larry Page, Founder, Google)
(Bill Gates, CollegeAttended, Harvard)

[Adapted example from Luke Zettlemoyer]

Distant Supervision [Mintz et al., ACL'09]

Corpus Text

Bill Gates founded Microsoft in 1975.

Bill Gates, founder of Microsoft, ...

Bill Gates attended Harvard from ...

Google was founded by Larry Page ...

Training Data

(Bill Gates, Microsoft)

Label: Founder

Feature: X founded Y

Freebase

(Bill Gates, Founder, Microsoft)

(Larry Page, Founder, Google)

(Bill Gates, CollegeAttended, Harvard)

[Adapted example from Luke Zettlemoyer]

Distant Supervision [Mintz et al., ACL'09]

Corpus Text

Bill Gates founded Microsoft in 1975.
Bill Gates, founder of Microsoft, ...
Bill Gates attended Harvard from ...
Google was founded by Larry Page ...

Training Data

(Bill Gates, Microsoft)
Label: Founder
Feature: X founded Y
Feature: X, founder of Y

Freebase

(Bill Gates, Founder, Microsoft)
(Larry Page, Founder, Google)
(Bill Gates, CollegeAttended, Harvard)

[Adapted example from Luke Zettlemoyer]

Distant Supervision [Mintz et al., ACL'09]

Corpus Text

Bill Gates founded Microsoft in 1975.
Bill Gates, founder of Microsoft, ...
Bill Gates attended Harvard from ...
Google was founded by Larry Page ...

Freebase

(Bill Gates, Founder, Microsoft)
(Larry Page, Founder, Google)
(Bill Gates, CollegeAttended, Harvard)

Training Data

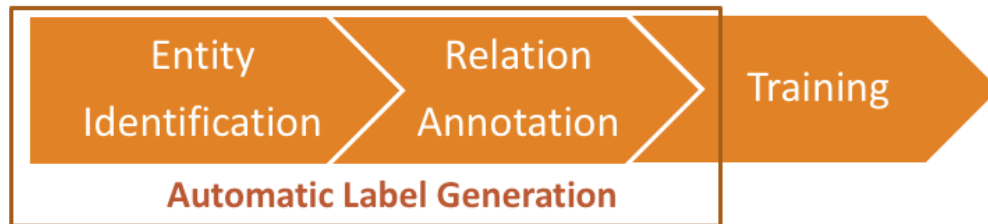
(Bill Gates, Microsoft)
Label: Founder
Feature: X founded Y
Feature: X, founder of Y

(Bill Gates, Harvard)
Label: CollegeAttended
Feature: X attended Y

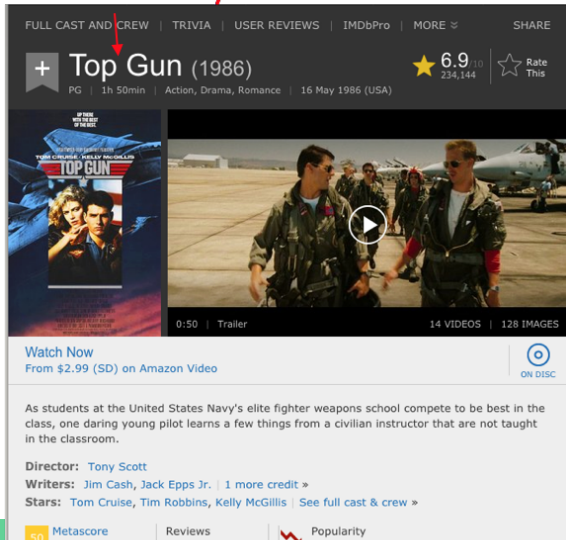
For negative examples, sample
unrelated pairs of entities.

[Adapted example from Luke Zettlemoyer]

Distantly Supervised Extraction--Ceres [Lockard et al., VLDB'18]

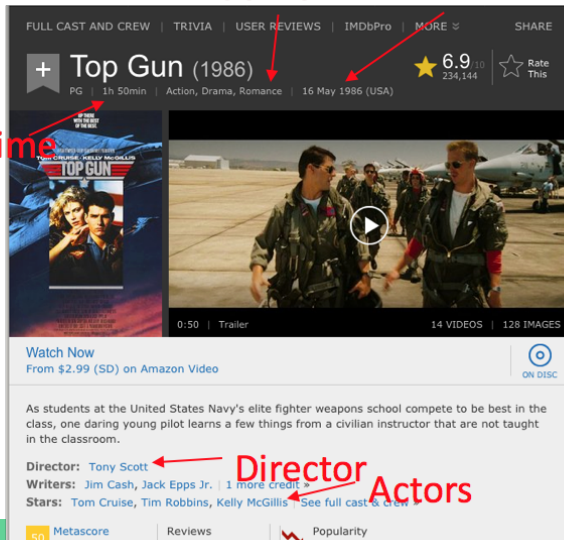


Movie entity



Genre Release Date

Runtime



Extracted triples

- (Top Gun, type.object.name, "Top Gun")
- (Top Gun, film.film.genre, Action)
- (Top Gun, film.film.directed_by, Tony Scott)
- (Top Gun, film.film.starring, Tom Cruise)
- (Top Gun, film.film.runtime, "1h 50min")
- (Top Gun, film.film.release_date_s, "16 May 1986")

Distantly Supervised Extraction--Ceres [Lockard et al., VLDB'18]

- Extraction experiments on SWDE benchmark

Vertical	Predicate	Vertex++			CERES-Full		
		P	R	F1	P	R	F1
Movie	Title	1.00	1.00	1.00	1.00	1.00	1.00
	Director	0.99	0.99	0.99	0.99	0.99	0.99
	Genre	0.88	0.87	0.87	0.93	0.97	0.95
	MPAA Rating	1.00	1.00	1.00	NA	NA	NA
	Average	0.97	0.97	0.97	0.97	0.99	0.98
NBAPlayer	Name	0.99	0.99	0.99	1.00	1.00	1.00
	Team	1.00	1.00	1.00	0.91	1.00	0.95
	Weight	1.00	1.00	1.00	1.00	1.00	1.00
	Height	1.00	1.00	1.00	1.00	0.90	0.95
	Average	1.00	1.00	1.00	0.98	0.98	0.98

Very high precision

Vertical	Predicate	Vertex++			CERES-Full		
		P	R	F1	P	R	F1
University	Name	1.00	1.00	1.00	1.00	1.00	1.00
	Type	1.00	1.00	1.00	0.72	0.80	0.76
	Phone	0.97	0.92	0.94	0.85	0.95	0.90
	Website	1.00	1.00	1.00	0.90	1.00	0.95
	Average	0.99	0.98	0.99	0.87	0.94	0.90
Book	Title	0.99	0.99	0.99	1.00	0.90	0.95
	Author	0.97	0.96	0.96	0.72	0.88	0.79
	Publisher	0.85	0.85	0.85	0.97	0.77	0.86
	Publication Date	0.90	0.90	0.90	1.00	0.40	0.57
	ISBN-13	0.94	0.94	0.94	0.99	0.19	0.32
	Average	0.93	0.93	0.93	0.94	0.63	0.70

Competent w. Wrapper induction w. manual annotation

Distantly Supervised Extraction--Ceres [Lockard et al., VLDB'18]

- Extraction on long-tail movie websites

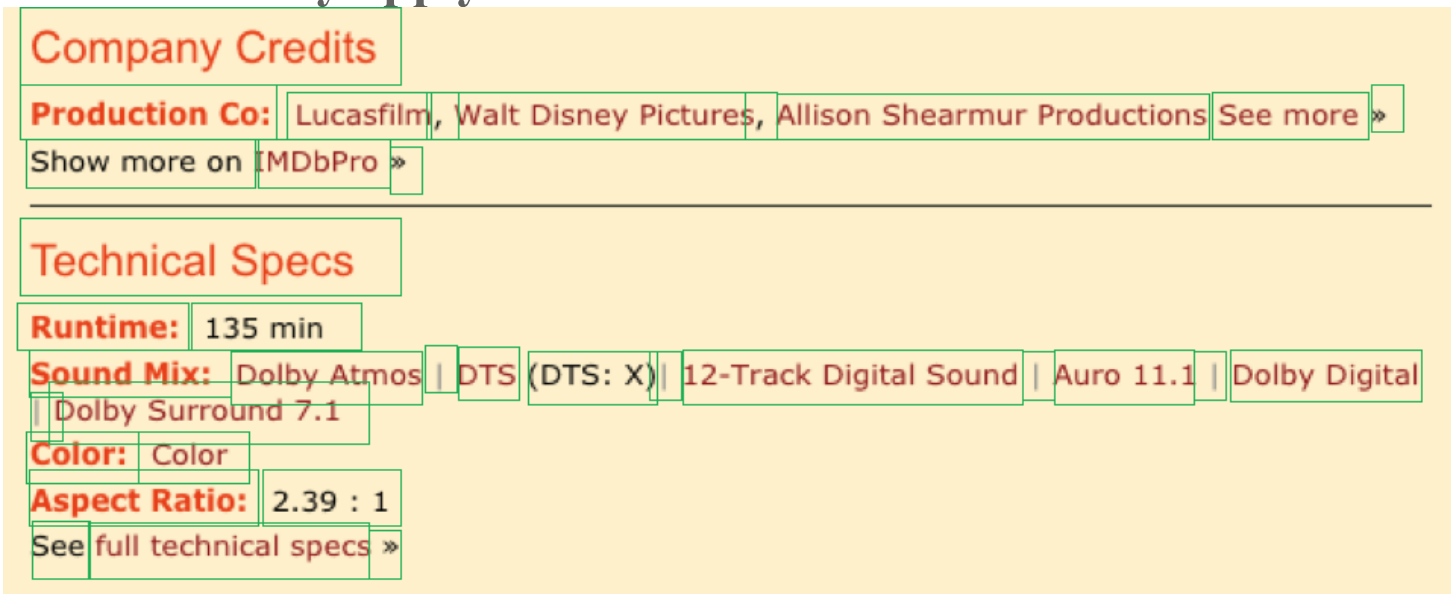
#Websites / #Webpages	33 / 434K
Language	English and 6 other languages
Domains	Animated films, Documentary films, Financial performance, etc.
# Annotated pages	70K (16%)
Annotated : Extracted #entities	1 : 2.6
Annotated : Extracted #triples	1 : 3.0
# Extractions	1.25 M
Precision	90%

Distantly Supervised Extraction--Ceres [Lockard et al., VLDB'18]

- **Which model is the best?**
 - Logistic regression: best results (20K features on one website)
 - Random forest: lower precision and recall
 - Deep learning??

Challenges in Applying Deep Learning on Extracting Semi-structured Data

- Web layout is neither 1D sequence nor regular 2D grid, so CNN or RNN does not directly apply



The image shows a screenshot of a movie's credits and technical specifications page. The page is divided into two main sections: "Company Credits" and "Technical Specs". Each section contains various pieces of information, some of which are highlighted with green bounding boxes. The "Company Credits" section includes the production companies (Lucasfilm, Walt Disney Pictures, Allison Shearmur Productions) and a link to "See more". The "Technical Specs" section includes the runtime (135 min), sound mix (Dolby Atmos, DTS, 12-Track Digital Sound, Auro 11.1, Dolby Digital), color (Color), and aspect ratio (2.39 : 1). A link to "See full technical specs" is also present.

Company Credits

Production Co: Lucasfilm, Walt Disney Pictures, Allison Shearmur Productions [See more »](#)

[Show more on IMDbPro »](#)

Technical Specs

Runtime: 135 min

Sound Mix: Dolby Atmos | DTS (DTS: X) | 12-Track Digital Sound | Auro 11.1 | Dolby Digital
| Dolby Surround 7.1

Color: Color

Aspect Ratio: 2.39 : 1

[See full technical specs »](#)



snorkel

Example System: Fonduer [Wu et al., SIGMOD'18]

Transistor Datasheet

SMBT3904...MMBT3904

NPN Silicon Switching Transistors

- High DC current gain: 0.1 mA to 100 mA
- Low collector-emitter saturation voltage

Maximum Ratings

Parameter	Symbol	Value	Unit
Collector-emitter voltage	V_{CE0}	40	V
Collector-base voltage	V_{CBO}	60	
Emitter-base voltage	V_{EBO}	6	
Collector current	I_C	200	mA
Total power dissipation	P_{tot}	330	mW
		250	
Temperature	T_i	150	°C
Storage temperature	T_{stg}	-65 ... 150	

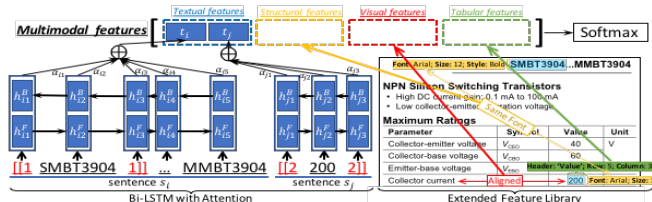
HasCollectorCurrent
(Transistor Part, Current)

SMBT3904 200mA

MMBT3904 200mA



Richly formatted data: information are expressed via textual, structural, tabular, and visual cues.



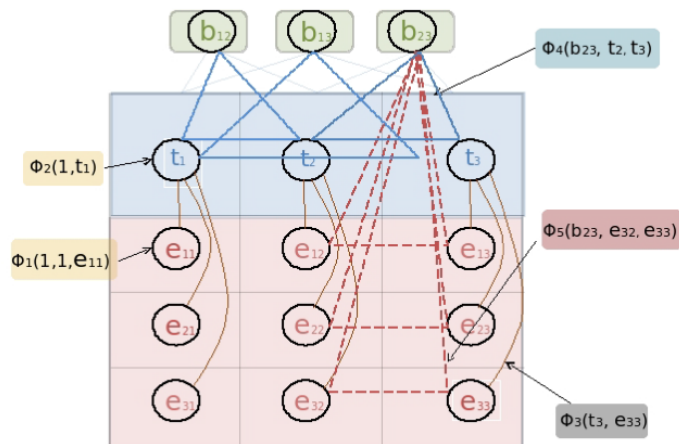
Fonduer combines a new bi-directional LSTM with multimodal features and weak supervision (specifically data programming).

Attend the talk in Research Session 13!

New version of code coming soon: <https://github.com/HazyResearch/fonduer>

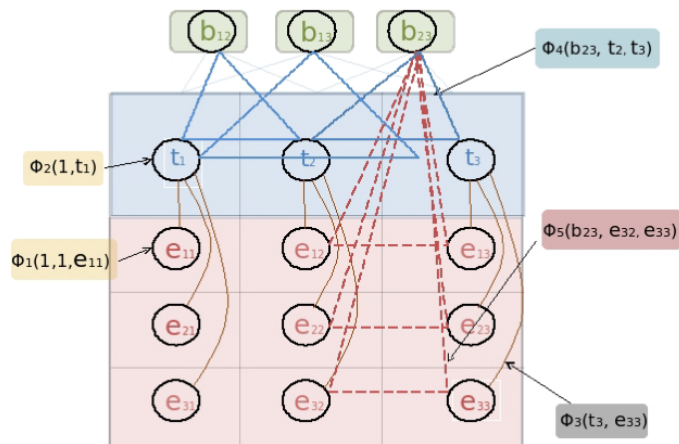
WebTable Extraction [Limaye et al., VLDB'10]

- Model table annotation using interrelated random variables, represented by a probabilistic graphical model
 - Cell text (in Web table) and entity label (in catalog)
 - Column header (in Web table) and type label (in catalog)
 - Column type and cell entity (in Web table)



WebTable Extraction [Limaye et al., VLDB'10]

- Model table annotation using interrelated random variables, represented by a probabilistic graphical model
 - Pair of column types (in Web table) and relation (in catalog)
 - Entity pairs (in Web table) and relation (in catalog)

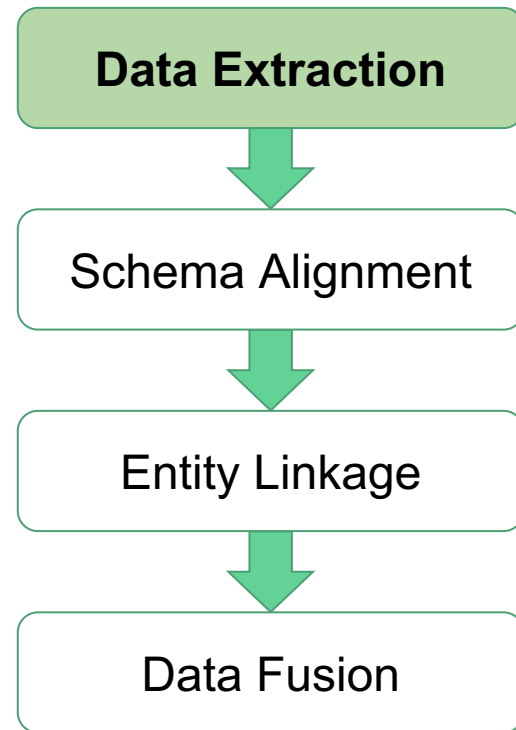


Challenges in Applying ML on DX

- Automatic data extraction cannot reach production quality requirement. How to improve precision?
- Every web designer has her own whim, but there are underlying patterns across websites. How to learn extraction patterns on different websites, especially for semi-structured sources?
- ClosedIE throws away too much data. How to apply OpenIE on all kinds of data?

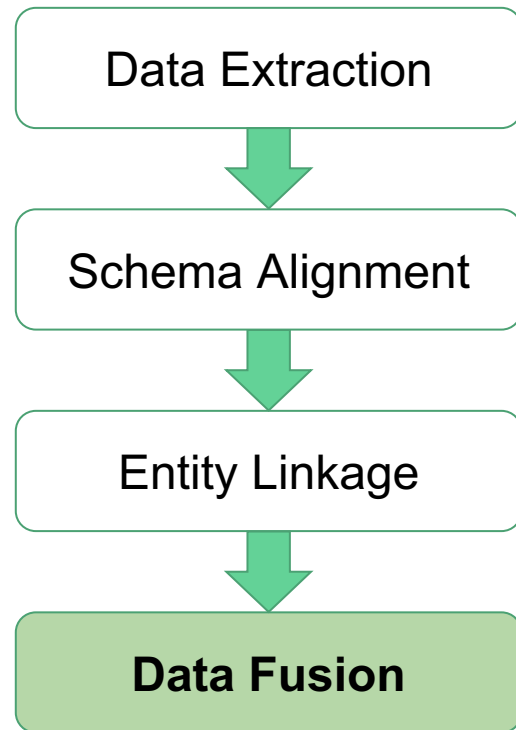
Recipe for Data Extraction

- **Problem definition: Extract structure from semi- or un-structured data**
- **Short answers**
 - **Wrapper induction** has high prec/rec
 - **Distant supervision is critical for collecting training data**
 - **LR is often effective; more research is needed for DL**



Outline

- Part I. Introduction
- Part II. ML for DI
 - ML for entity linkage
 - ML for data extraction
 - ML for data fusion
 - ML for schema alignment
- Part III. DI for ML
- Part IV. Conclusions and research direction




What is Data Fusion?


- **Definition:** Resolving conflicting data and verifying facts.
- **Example:** “*OK Google, How long is the Mississippi River?*”

Mississippi River / Length

2,320 mi

People also search for

 Missouri River
2.341K mi

 Nile
4.258K mi

Mississippi River

River in the United States of America

4.2 ★★★★★ 400 Google reviews

The Mississippi River is the chief river of the second-largest drainage system on the North American continent, second only to the Hudson Bay drainage system.

[Wikipedia](#)

Discharge: 593,000 cubic feet per second

Basin area: 1.151 million mi²

Source: [Lake Itasca](#)

Mouth: [Gulf of Mexico](#)

Country: [United States of America](#)

Did you know: The Mississippi River is the second-longest river in the US (2,202 mi). [wikipedia.org](#)

Mississippi River Facts - Mississippi National River and Recreation ...

<https://www.nps.gov/miss/riverfacts.htm>

Nov 14, 2017 - The staff of Itasca State Park at the Mississippi's headwaters suggest the main stem of the river is 2,552 miles long. The US Geologic Survey has published a number of 2,300 miles, the EPA says it is 2,320 miles long, and the Mississippi National River and Recreation Area suggests the river's length is 2,350 miles.

Longest mainstem rivers of the United States									
#	Name	Mouth ^[5]	Length	Source coordinates ^[11]	Mouth coordinates ^[11]	Watershed area ^[12]	Discharge ^[12]	States, provinces, and image ^{[8][11]}	
1	Missouri River	Mississippi River	2,341 mi 3,768 km ^[13]	 45°55'39"N 111°30'29"W ^[14]	 38°48'49"N 90°07'11"W	529,353 mi ² 1,371,017 km ² ^[15] ↓ ^[n 2]	69,100 ft ³ /s 1,956 m ³ /s [n 3]	Montana ^a , North Dakota, South Dakota, Nebraska, Iowa, Kansas, Missouri ^m	
2	Mississippi River	Gulf of Mexico	2,202 mi 3,544 km ^[17] [n 4]	 47°14'22"N 95°12'29"W ^[18]	 29°09'04"N 89°15'12"W	1,260,000 mi ² 3,270,000 km ² ^[19] ↓ ^[n 5]	650,000 ft ³ /s 18,400 m ³ /s	Minnesota ^a , Wisconsin, Iowa, Illinois, Missouri, Kentucky, Tennessee, Arkansas, Mississippi, Louisiana ^m	

The Basic Setup of Data Fusion

Source Observations

Source	River	Attribute	Value
KG	Mississippi River	Length	2,320 mi
KG	Missouri River	Length	2,341 mi
Wikipedia	Mississippi River	Length	2,202 mi
Wikipedia	Missouri River	Length	2,341 mi
USGS	Mississippi River	Length	2,340 mi
USGS	Missouri River	Length	2,540 mi

Fact

Conflicting value

Source reports
a value for a fact

True Facts

River	Attribute	Value
Mississippi River	Length	?
Missouri River	Length	?

Fact's true value

Goal: Find the **latent**
true value of facts.

The Basic Setup of Data Fusion

Source Observations

Source	River	Attribute	Value
KG	Mississippi River	Length	2,320 mi
KG	Missouri River	Length	2,341 mi
Wikipedia	Mississippi River	Length	2,202 mi
Wikipedia	Missouri River	Length	2,341 mi
USGS	Mississippi River	Length	2,340 mi
USGS	Missouri River	Length	2,540 mi

Fact

Conflicting value

Source reports
a value for a fact

True Facts

River	Attribute	Value
Mississippi River	Length	?
Missouri River	Length	?

Fact's true value

Idea: Use *redundancy* to infer the true value of each fact.

Majority Voting for Data Fusion

Source Observations

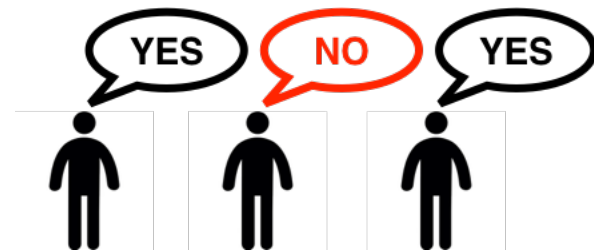
Source	River	Attribute	Value
KG	Mississippi River	Length	2,320 mi
KG	Missouri River	Length	2,341 mi
Wikipedia	Mississippi River	Length	2,202 mi
Wikipedia	Missouri River	Length	2,341 mi
USGS	Mississippi River	Length	2,340 mi
USGS	Missouri River	Length	2,540 mi

Majority voting can be limited. What if sources are correlated (e.g., copying)?

Idea: Model source quality for accurate results.

True Facts

River	Attribute	Value
Mississippi River	Length	?
Missouri River	Length	2,341



MV's assumptions

1. Sources report values independently
2. Sources are better than chance.

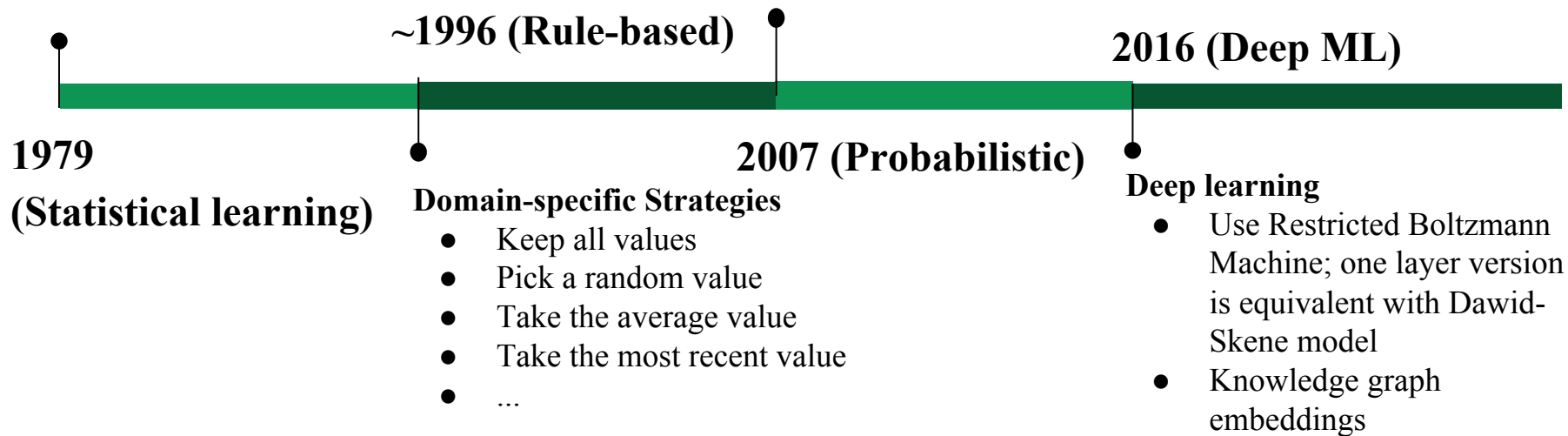
40 Years of Data Fusion (beyond Majority Voting)

Dawid-Skene model

- Model the error-rate of sources
- Expectation-maximization

Probabilistic Graphical Models

- Use of generative models
- Focus on unsupervised learning



A Probabilistic Model for Data Fusion

- **Random variables:** Introduce a *latent random variable* to represent the true value of each fact.
- **Features:** Source observations become features associated with different random variables.
- **Model parameters:** Weights related to the error-rates of each data source.

$$P(\text{Fact} = v | \text{data}) = \underbrace{\frac{1}{Z}}_{\text{Normalizing constant}} \exp \sum_{s \in \text{Sources}} \sum_{v' \in \text{Values}} \sigma_S^{v,v'} \cdot 1[S \text{ reports Fact} = v']$$

error-rate scores (model parameters)

$$\sigma_S^{v,v'} = \log \left(\frac{\text{Error-rate of Source } S}{1 - \text{Error-rate of Source } S} \right)$$

Error-rate = probability that a source provides value v' instead of value v

The Challenge of Training Data

- How much data do we need to train the data fusion model?
- **Theorem:** We need a number of labeled examples proportional to the number of sources [Ng and Jordan, NIPS'01]
- **Model parameters:** Weights related to the error-rates of each data source.

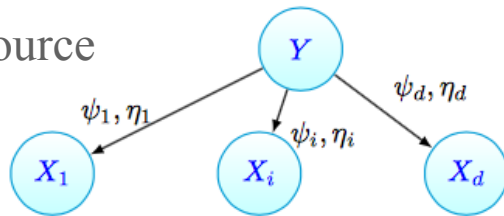
But the number of sources can be in the thousands or millions
and training data is limited!

Idea 1: Leverage redundancy and use unsupervised learning.

The Dawid-Skene Algorithm [Dawid and Skene, 1979]

Iterative process to estimate data source error rates

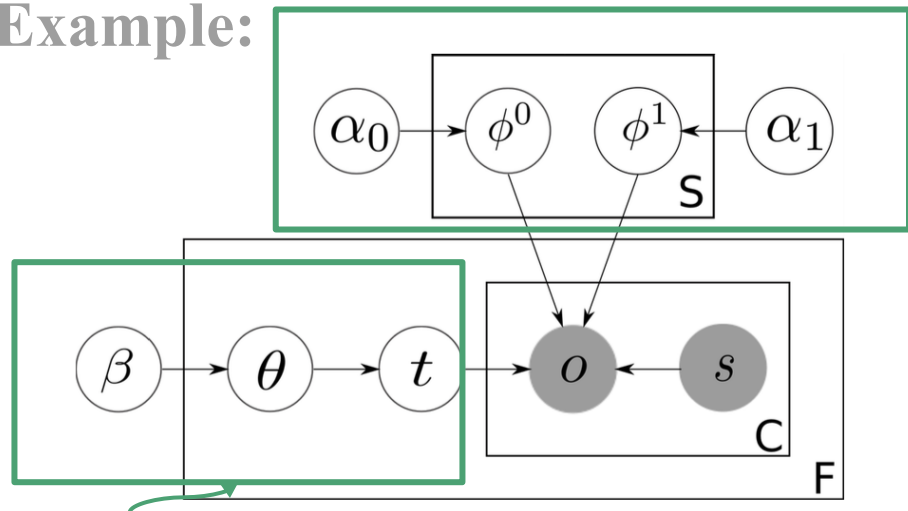
1. Initialize “inferred” true value for each fact (e.g., use majority vote)
2. Estimate **error rates** for workers (using “inferred” true values)
3. Estimate **“inferred” true values** (using error rates, weight source votes according to quality)
4. Go to Step 2 and iterate until convergence



Assumptions: (1) average source error rate < 0.5 , (2) dense source observations, (3) conditional independence of sources, (4) errors are uniformly distributed across all instances.

Probabilistic Graphical Models for Data Fusion

Example:



Prior truth probability [Zhao et al., VLDB 2012]

Source Quality

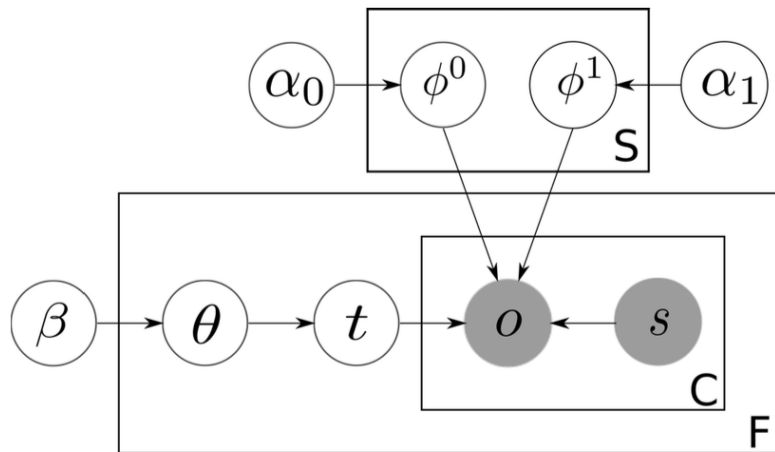
Setup: Identify true source claims

Entity (Movie)	Attribute (Cast)	Source
Harry Potter	Daniel Radcliffe	IMDB
Harry Potter	Emma Waston	IMDB
Harry Potter	Rupert Grint	IMDB
Harry Potter	Daniel Radcliffe	Netflix
Harry Potter	Daniel Radcliffe	BadSource.com
Harry Potter	Emma Waston	BadSource.com
Harry Potter	Johnny Depp	BadSource.com
Pirates 4	Johnny Depp	Hulu.com
...

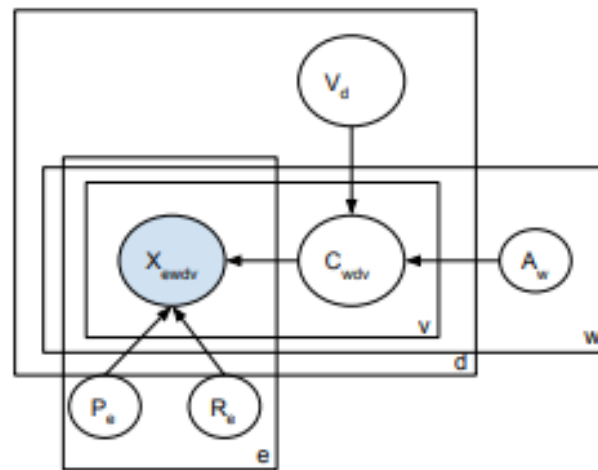
Extensive work on modeling source observations and source interactions to address limitations of basic Dawid-Skene.

Probabilistic Graphical Models for Data Fusion

Modeling both source quality and extractor accuracy



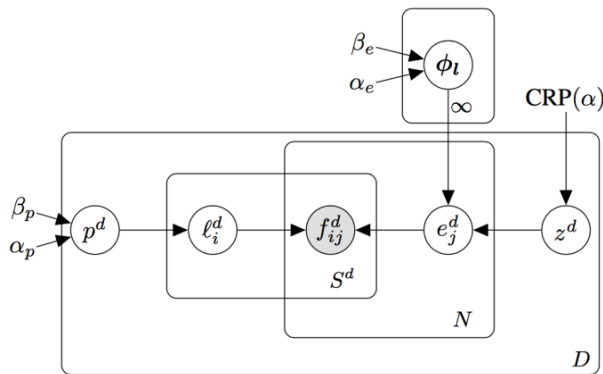
[Zhao et al., VLDB 2012]



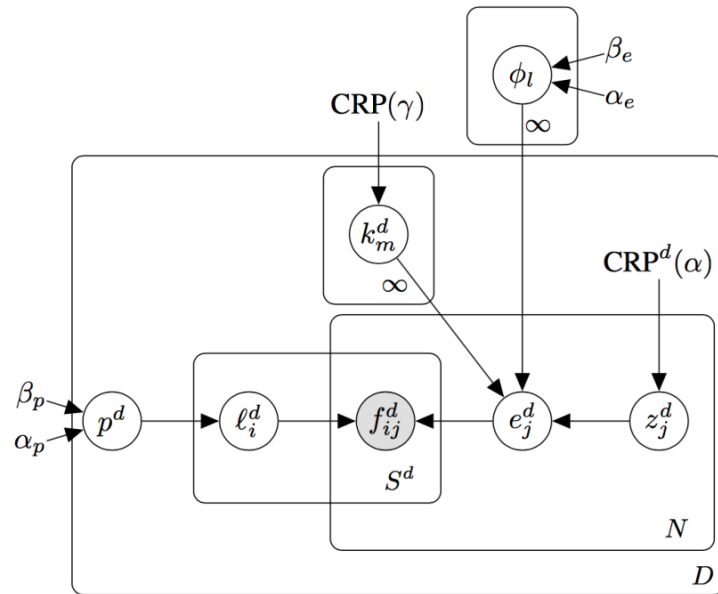
[Dong et al., VLDB 2015]

Extensive work on modeling source observations and source interactions to address limitations of basic Dawid-Skene.

Probabilistic Graphical Models for Data Fusion



Modeling source dependencies



[Platanios et al., ICML 2016]

Extensive work on modeling source observations and source interactions to address limitations of basic Dawid-Skene.

PGMs in Data Fusion [Li et al., VLDB'14]

Table 6: Summary of data-fusion methods. X indicates that the method considers the particular evidence.

Category	Method	#Providers	Source trustworthiness	Item trustworthiness	Value Popularity	Value similarity	Value formatting	Copying
Baseline	Vote	X						
Web-link based	HUB	X	X					
	AVGLOG	X	X					
	INVEST	X	X					
	POOLEDINVEST	X	X					
IR based	2-ESTIMATES	X	X					
	3-ESTIMATES	X	X	X				
	COSINE	X	X					
Bayesian based	TRUTHFINDER	X	X			X		
	ACCUPR	X	X					
	POPACCU	X	X		X			
	ACCUSIM	X	X			X		
	ACCUFORMAT	X	X			X	X	
Copying affected	ACCUCOPY	X	X			X	X	X

Bayesian models capture source observations and source interactions.

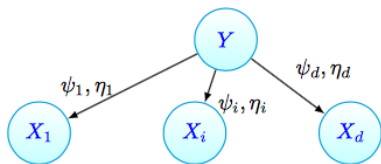
PGMs in Data Fusion [Li et al., VLDB'14]

Category	Method	<i>Stock</i>				<i>Flight</i>			
		prec w. trust	prec w/o. trust	Trust dev	Trust diff	prec w. trust	prec w/o. trust	Trust dev	Trust diff
Baseline	Vote	-	.908	-	-	-	.864	-	-
Web-link based	HUB	.913	.907	.11	.08	.939	.857	.2	.14
	AVGLOG	.910	.899	.17	-.13	.919	.839	.24	.001
	INVEST	.924	.764	.39	-.31	.945	.754	.29	-.12
	POOLEDINVEST	.924	.856	1.29	0.29	.945	.921	17.26	7.45
IR based	2-ESTIMATES	.910	.903	.15	-.14	.87	.754	.46	-.35
	3-ESTIMATES	.910	.905	.16	-.15	.87	.708	.95	-.94
	COSINE	.910	.900	.21	-.17	.87	.791	.48	-.41
Bayesian based	TRUTHFINDER	.923	.911	.15	.12	.957	.793	.25	.16
	ACCUPr	.910	.899	.14	-.11	.91	.868	.16	-.06
	POPACCU	.909	.892	.14	-.11	.958	.925	.17	-.11
	ACCUSIM	.918	.913	.17	-.16	.903	.844	.2	-.09
	ACCUFORMAT	.918	.911	.17	-.16	.903	.844	.2	-.09
	ACCUSIMATTR	.950	.929	.17	-.16	.952	.833	.19	-.08
	ACCUFORMATATTR	.948	.930	.17	-.16	.952	.833	.19	-.08
Copying affected	ACCUCOPY	.958	.892	.28	-.11	.960	.943	.16	-.14

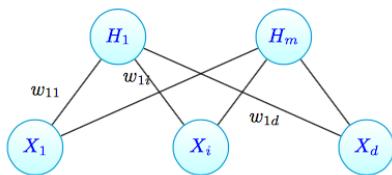
Modeling the quality of data sources leads to improved accuracy.

Dawid-Skene and Deep Learning [Shaham et al., ICML'16]

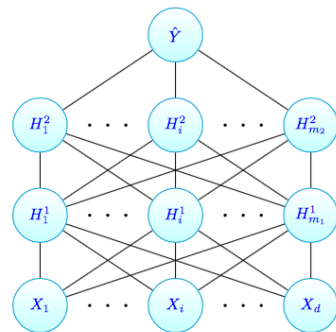
Theorem: The Dawid and Skene model is *equivalent* to a Restricted Boltzmann Machine (RBM) with a single hidden node.



Dawid and Skene model.



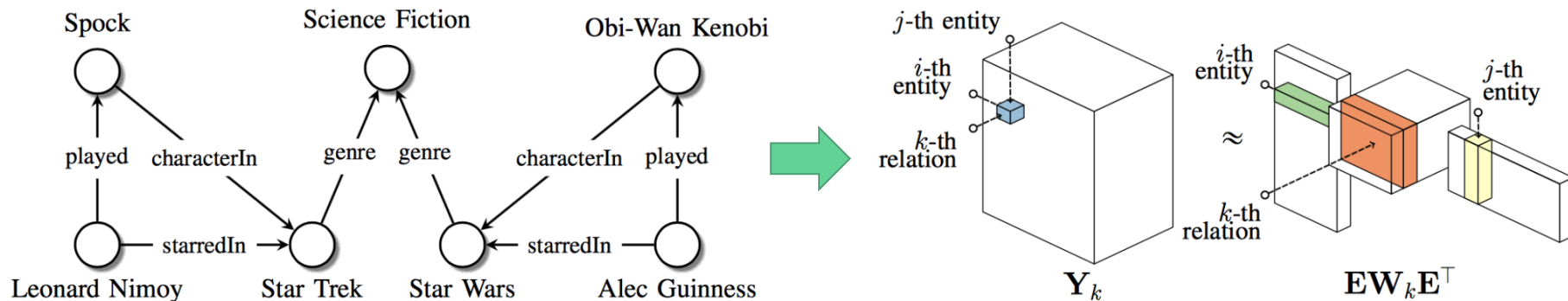
A RBM with d visible and m hidden units.



Sketch of a two-hidden-layer RBM-based DNN.

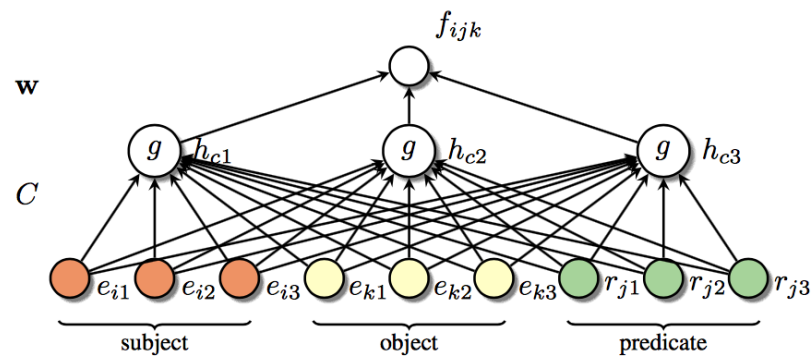
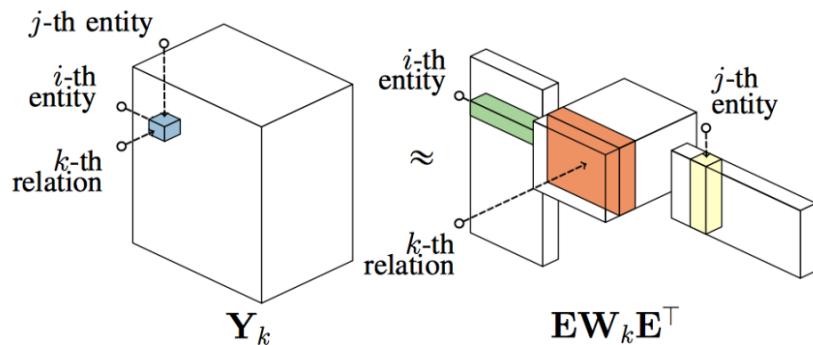
When the conditional independence assumption of Dawid-Skene does not hold, a better approximation may be obtained from a deeper network.

Knowledge Graph Embeddings [Survey: Nickiet et al., 2015]



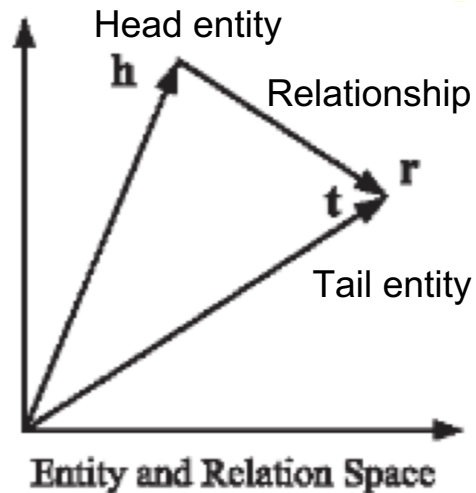
A knowledge graph can be encoded as a tensor.

Knowledge Graph Embeddings [Survey: Nickiet et al., 2015]



Neural networks can be used to obtain richer representations.

Knowledge Graph Embeddings



Example: Learn embeddings from IMDB data and identify various types of errors in WikiData [Dong et al., KDD'18]

Subject	Relation	Target	Reason
The Moises Padilla Story	writtenBy	César Ámigo Aguilar	Linkage error
Bajrangi Bhaijaan	writtenBy	Yo Yo Honey Singh	Wrong relationship
Piste noire	writtenBy	Jalil Naciri	Wrong relationship
Enter the Ninja	musicComposedBy	Michael Lewis	Linkage error
The Secret Life of Words	musicComposedBy	Hal Hartley	Cannot confirm
...

- TransE: $\text{score}(h,r,t) = -\|h+r-t\|_{1/2}$
- Hot field with increasing interest [Survey by Wang et al., TKDE 2017]

The Challenge of Training Data

- How much data do we need to train the data fusion model?
- **Theorem:** We need a number of labeled examples proportional to the number of sources [Ng and Jordan, NIPS'01]
- **Model parameters:** Weights related to the error-rates of each data source.

But the number of sources can be in the thousands or millions
and training data is limited!

Idea 1: Leverage redundancy and used unsupervised learning.

Idea 2: Limit model parameters and use a small number of training data.

SLiMFast: Discriminative Data Fusion [Rekatsinas et al., SIGMOD'17]

Limit the informative parameters of the model by using domain knowledge

Key Idea: Sources have (domain specific) features that are indicative of error rates

Example:

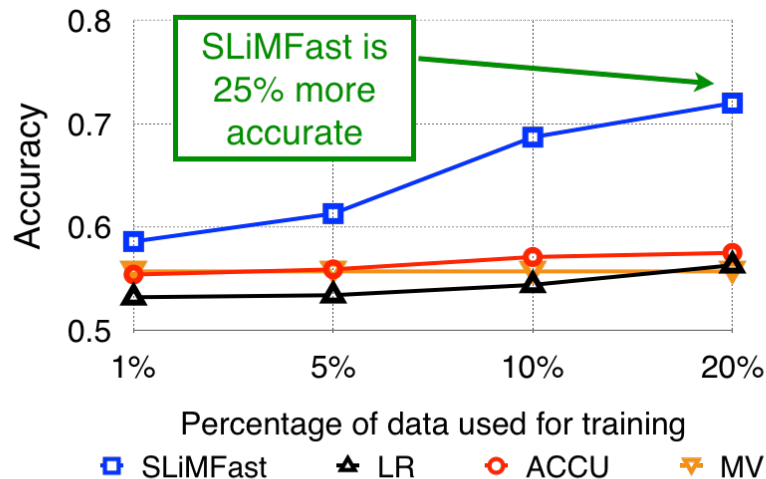
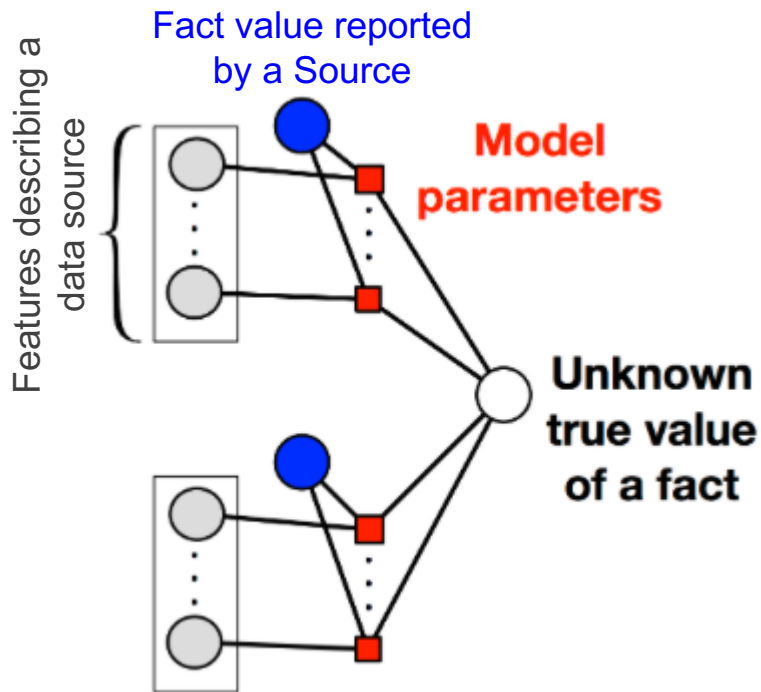


- newly registered similar to existing domain
- traffic statistics
- text quality (e.g., misspelled words, grammatical errors)
- sentiment analysis



- avg. time per task
- number of tasks
- market used

SLiMFast: Discriminative Data Fusion [Rekatsinas et al., SIGMOD'17]



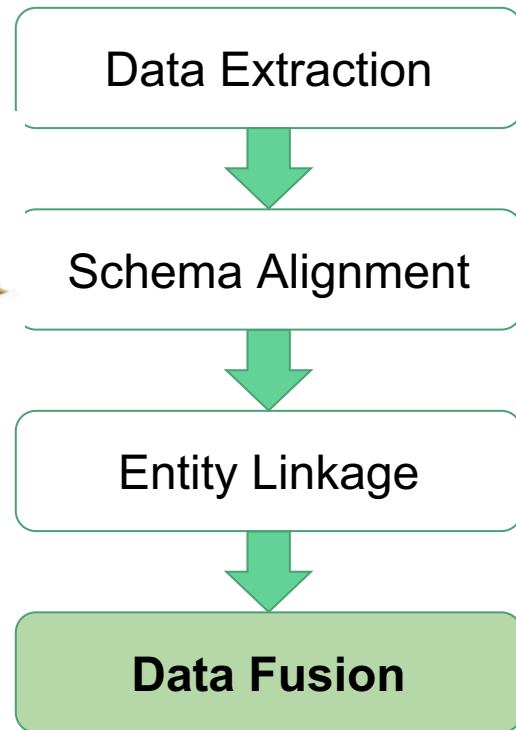
Genomics data: 2.7k sources (articles), 571 objects (gene-disease), 4 domain features (year, citation, author, journal)

Challenges in Data Fusion

- There are few solutions for unstructured data. Mostly work on fact verification [Tutorial by Dong et al., KDD'2018]. Most data Fusion solutions assume data extraction. Can state-of-the art DL help?
- Using training data is key and semi-supervised learning can significantly improve the quality of Data Fusion results. How can one collect training data effectively without manual annotation?
- We have only scratched the surface of what representation learning and deep learning methods can offer. Can deep learning streamline data fusion? What are its limitations?

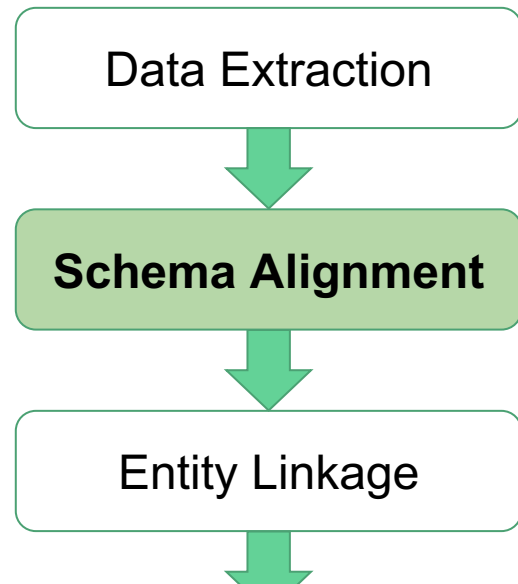
Recipe for Data Fusion

- **Problem definition: Resolve conflicts and obtain correct values**
- **Short answers**
 - Reasoning about source quality is key and works for easy cases
 - Semi-supervised learning has shown **BIG** potential
 - Representation learning provides positive evidence for streamlining data fusion.



Outline

- Part I. Introduction
- Part II. ML for DI
 - ML for entity linkage
 - ML for data extraction
 - ML for data fusion
 - ML for schema alignment
- Part III. DI for ML



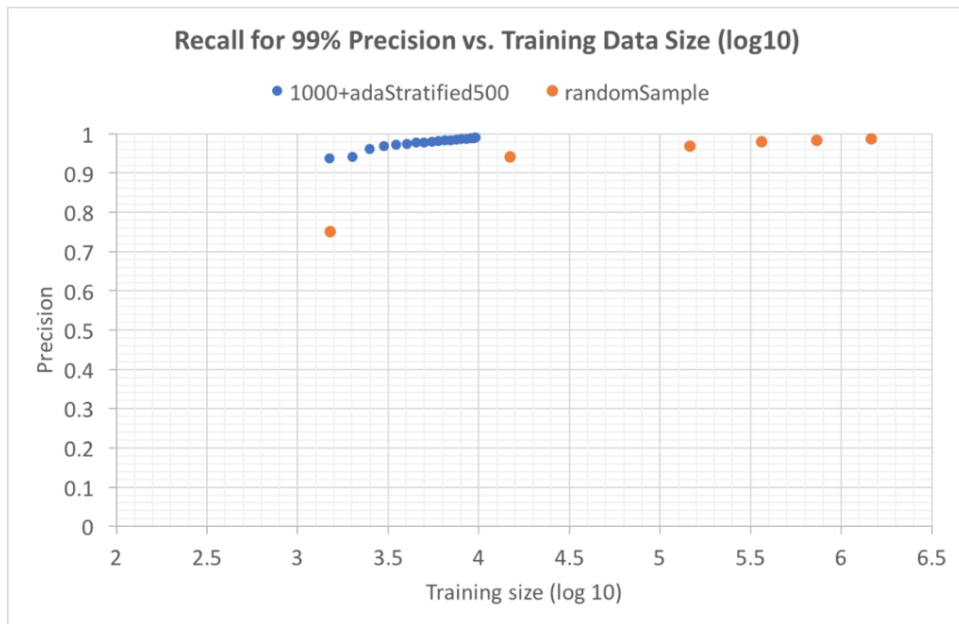
Come to our VLDB tutorial for schema alignment and universal schema!!

Revisit Theme I. Which ML Model Works Best?

DI tasks	Hyperplanes (e.g., Log Reg)	Kernal (e.g., SVM)	Tree-based (e.g., Random forest)	Graphical models (e.g., CRF)	Logic programs (e.g, soft logic)	Neural networks (e.g., RNN)
Entity resolution	X	X	X		X	X
Data fusion	X			X		
DOM extraction	X					
Text extraction	X			X		X
Schema alignment	X			X		X

No single winner, although ensemble models and deep learning models show promising results.

Revisit Theme II. Does Supervised Learning Apply to DI?

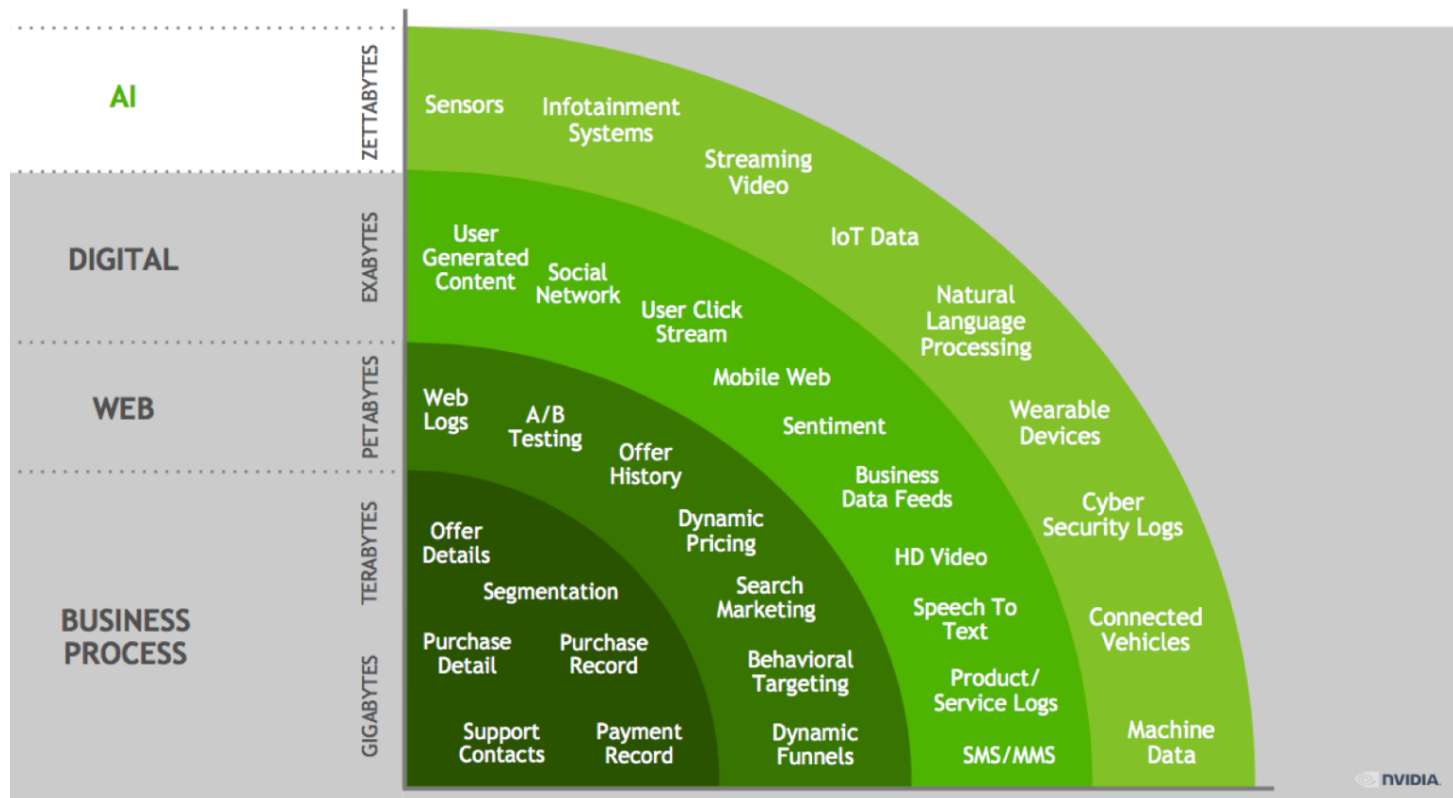


Active learning, semi-supervised methods, and weak supervision lead to dramatically more efficient solutions.

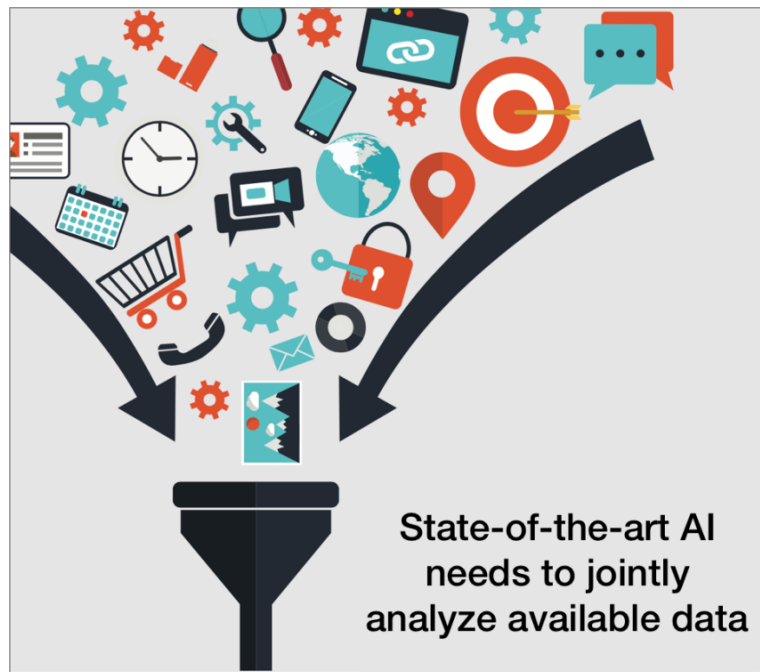
Outline

- Part I. Introduction
- Part II. ML for DI
- **Part III. DI for ML**
 - Training data creation
 - Data cleaning
- Part IV. Conclusions and research directions

ML is data-hungry



Successful ML requires Data Integration



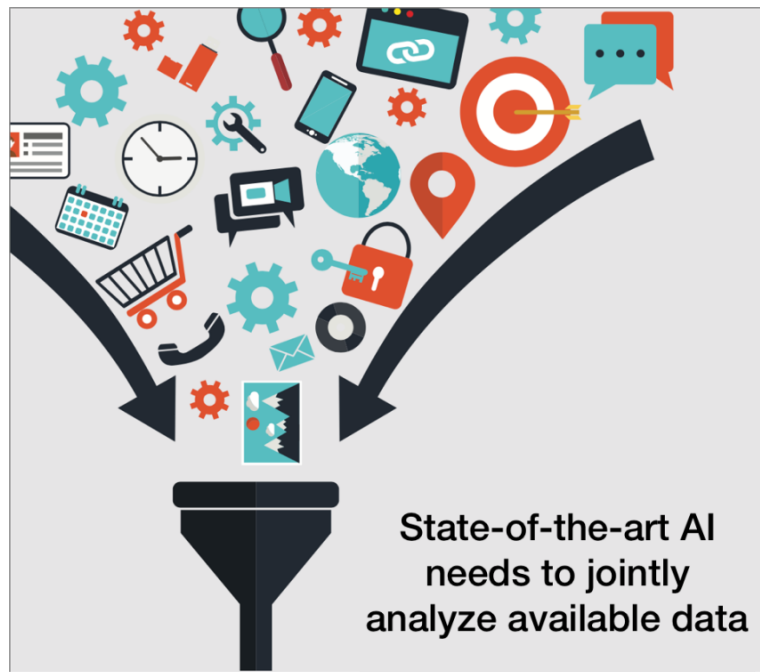
IMAGENET MovieLens



COCO is a large-scale object detection, segmentation, and captioning dataset.

Large collections of manually curated training data are necessary for progress in ML.

Successful ML requires Data Integration



IMAGENET MovieLens



COCO is a large-scale object detection, segmentation, and captioning dataset.

Large collections of manually curated **training data** are necessary for progress in ML.

Outline

- Part I. Introduction
- Part II. ML for DI
- Part III. DI for ML
 - Training data creation
 - Data cleaning
- Part IV. Conclusions and research directions

50 Years of Artificial Intelligence

Expert systems

- Manually curated knowledge bases of facts and rules
- Use of inference engines
- No support for high-dimensional data

Graphical models and logic

- Relational statistical learning
- Markov logic network

2010s

(Representation Learning)

1990s (Features)

2009 (PGMs)

1970s (Rules)

Classical ML

- Low complexity models
- Strong priors that capture domain knowledge (feature engineering)
- Small amounts of training data

Deep learning

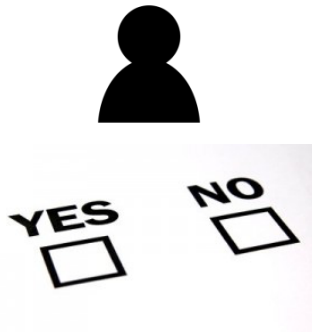
- Automatically learn representations
- Impressive with high-dimensional data
- Data hungry!

The ML Pipeline in the Deep Learning Era

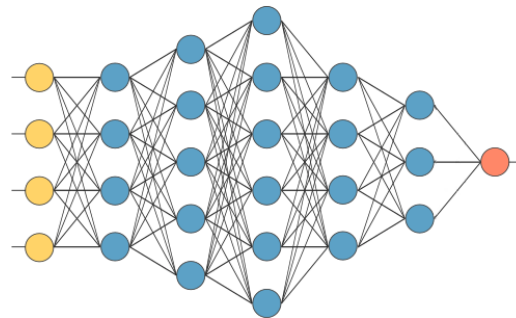
Data Collection



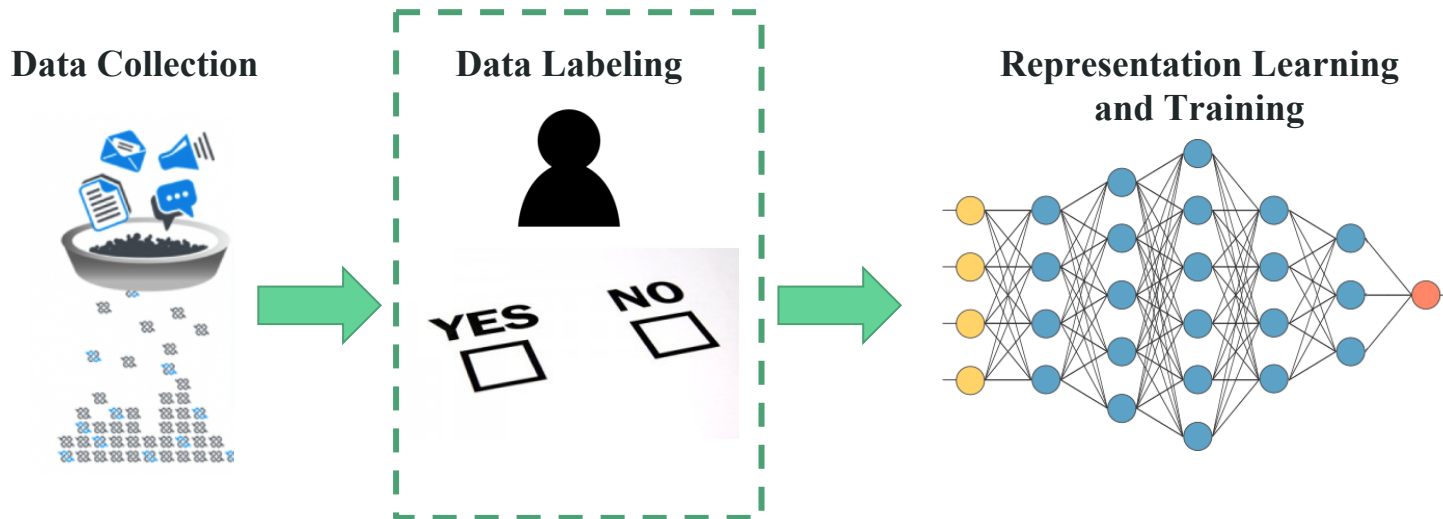
Data Labeling



**Representation Learning
and Training**



The ML Pipeline in the Deep Learning Era



Main pain point today, most time spent in labeling data.

Training Data: Challenges and Opportunities

- Collecting training data is **expensive** and **slow**.
- We are overfitting to our training data. [Recht et al., 2018]
 - Hand-labeled training data does not change
- Training data is the point to inject domain knowledge
 - Modern ML is too complex to hand-tune features and priors

Training Data: Challenges and Opportunities

- Collecting training data is **expensive** and **slow**.
- We are overfitting to our training data. [Recht et al., 2018]
 - Hand-labeled training data does not change
- Training data is the point to inject domain knowledge
 - Modern ML is too complex to hand-tune features and priors

How do we get training data more effectively?

The Rise of Weak Supervision

Definition: Supervision with noisy (much easier to collect) labels; prediction on a larger set, and then training of a model.

Semi-supervised learning and ensemble learning

Examples:

- use of non-expert labelers (crowdsourcing),
- use of curated catalogs (distant supervision)
- use of heuristic rules (labeling functions)



NELL



snorkel

The Rise of Weak Supervision

- Alexa – Customer embrace of Alexa continues, with Alexa-enabled devices among the best-selling items across all of Amazon. We're seeing extremely strong adoption by other companies and developers that want to create their own experiences with Alexa. There are now more than 30,000 skills for Alexa from outside developers, and customers can control more than 4,000 smart home devices from 1,200 unique brands with Alexa. The foundations of Alexa continue to get smarter every day too. We've developed and implemented an on-device fingerprinting technique, which keeps your device from waking up when it hears an Alexa commercial on TV. (This technology ensured that our Alexa Super Bowl commercial didn't wake up millions of devices.) Far-field speech recognition (already very good) has improved by 15% over the last year; and in the U.S., U.K., and Germany, we've improved Alexa's spoken language understanding by more than 2% over the last 12 months through enhancements in Alexa's machine learning components and the use of semi-supervised learning techniques. (These semi-supervised learning techniques reduced the amount of labeled data needed to achieve the same accuracy improvement by 40 times!) Finally, we've dramatically reduced the amount of time required to teach Alexa new languages by using machine translation and transfer learning techniques, which allows us to serve customers in more countries (like India and Japan).

The Rise of Weak Supervision

Definition: Supervision with noisy (much easier to collect) labels; prediction on a larger set, and then training of a model.

Related to semi-supervised learning and ensemble learning

Examples: use of non-expert labelers (crowdsourcing), use of curated catalogs (distant supervision), use of heuristic rules (labeling functions)

Methods developed to tackle data integration problems are closely related to weak supervision.

Learning from Crowds [Raykar et al., JMLR'10]

Setup: Supervised learning but instead of gold groundtruth one has access to multiple annotators providing (possibly noisy) labels (no absolute gold standard).

Task: Learn a classifier from multiple noisy labels.

Closely related to Dawid-Skene!

Difference: Estimating the ground truth and the annotator performance is a byproduct here. Goal is to learn a classifier.

Learning from Crowds [Raykar et al., JMLR'10]

Example Task: Binary classification

$$\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$$

N examples, with labels $\mathbf{y}_i = y_i^1, \dots, y_i^R$
provided by R different annotators

Learning from Crowds [Raykar et al., JMLR'10]

Example Task: Binary classification

Annotator performance:

Sensitivity (true positive rate)

$$\alpha^j = \Pr[y^j = 1 | y = 1]$$

$$\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$$

N examples, with labels $\mathbf{y}_i = y_i^1, \dots, y_i^R$
provided by R different annotators

Specificity (1 - false positive rate)

$$\beta^j = \Pr[y^j = 0 | y = 0]$$

Learning from Crowds [Raykar et al., JMLR'10]

Example Task: Binary classification

Annotator performance:

Sensitivity (true positive rate)

$$\alpha^j = \Pr[y^j = 1 | y = 1]$$

Specificity (1 - false positive rate)

$$\beta^j = \Pr[y^j = 0 | y = 0]$$

Learning: $\Pr[\mathcal{D} | \theta] = \prod_{i=1}^N [a_i p_i + b_i (1 - p_i)]$

$$p_i := \sigma(\mathbf{w}^\top \mathbf{x}_i).$$

$$a_i := \prod_{j=1}^R [\alpha^j]^{y_i^j} [1 - \alpha^j]^{1 - y_i^j}.$$

$$b_i := \prod_{j=1}^R [\beta^j]^{1 - y_i^j} [1 - \beta^j]^{y_i^j}.$$

Model
parameters
 $\{\mathbf{w}, \boldsymbol{\alpha}, \boldsymbol{\beta}\}$

EM algorithm to obtain maximum-likelihood estimates. Difference with Dawid-Skene is the estimation of \mathbf{w} .

Distant Supervision [Mintz et al., ACL'09]

Goal: Extracting structured knowledge from text.

Hypothesis: If two entities belong to a certain relation, any sentence containing those two entities is likely to express that relation.

Idea: Use a *database* of relations to gets lots of *noisy* training examples

- Instead of hand-creating seed tuples (bootstrapping)
- Instead of using hand-labeled corpus (supervised)

Benefits: has the advantages of supervised learning (leverage reliable hand-created knowledge), has the advantages of unsupervised learning (leverage unlimited amounts of text data).

Remember: Distant Supervision [Mintz et al., ACL'09]

Example task: Relation extraction.

Corpus Text

Bill Gates founded Microsoft in 1975.
Bill Gates, founder of Microsoft, ...
Bill Gates attended Harvard from ...
Google was founded by Larry Page ...

Freebase

Founder: (Bill Gates, Microsoft)
Founder: (Larry Page, Google)
CollegeAttended: (Bill Gates, Harvard)

Training Data

(Bill Gates, Microsoft)
Label: Founder
Feature: X founded Y
Feature: X, founder of Y

(Bill Gates, Harvard)
Label: CollegeAttended
Feature: X attended Y

For negative examples, sample
unrelated pairs of entities.

[Adapted example from Luke Zettlemoyer]

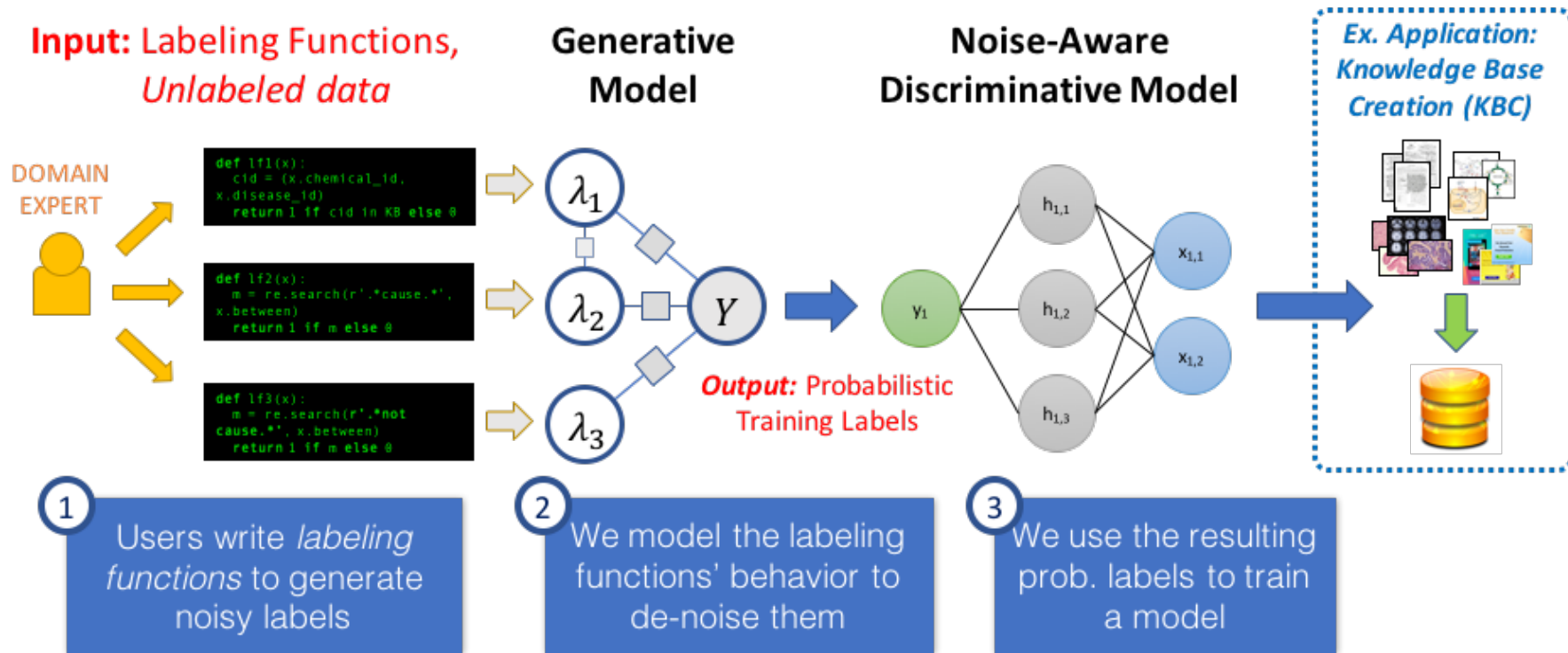
Distant Supervision [Mintz et al., ACL'09]

Entity Linking is an inherent problem in Distant Supervision.

The quality of matches can vary significantly and has a direct effect on extraction quality.

Relation	Freebase Matches	
	#sents	% true
/business/person/company	302	89.0
/people/person/place_lived	450	60.0
/location/location/contains	2793	51.0
/business/company/founders	95	48.4
/people/person/nationality	723	41.0
/location/neighborhood/neighborhood_of	68	39.7
/people/person/children	30	80.0
/people/deceased_person/place_of_death	68	22.1
/people/person/place_of_birth	162	12.0
/location/country/administrative_divisions	424	0.2

Snorkel: Code as Supervision [Ratner et al., NIPS'16, VLDB'18]



Snorkel: Code as Supervision [Ratner et al., NIPS'16, VLDB'18]

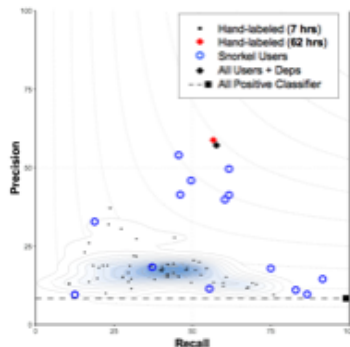


Snorkel biomedical workshop in collaboration with the NIH Mobilize Center



15 companies and research groups attended

How well did these new Snorkel users do?



71% New Snorkel users matched or beat 7 hours of hand-labeling

2.8x Faster than hand-labeling data

45.5% Average improvement in model performance



3rd Place Score

No machine learning experience
Beginner-level Python

[Slide by Alex Ratner]

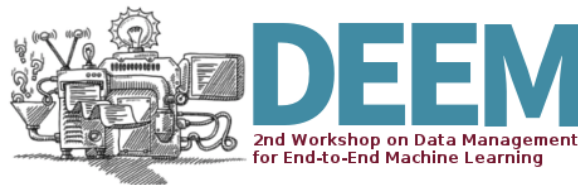
Alex (the creator of Snorkel) is on the market!

Alex Ratner



<https://ajratner.github.io>

Find out more about Snorkel
MeTaL and weak supervision
for Multi-task Learning at



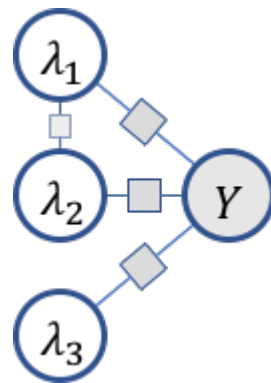
Friday in Montgomery

Challenges in Creating Training Data

- Richly-formatted data is still a challenge. How can attack weak supervision when data includes images, text, tables, video, etc.?
- Combining weak supervision with other data enrichment techniques such as data augmentation is an exciting direction. How can reinforcement learning help here (<http://goo.gl/K2qopQ>)?
- How can we combine weak supervision with techniques from semi-supervised?
- Most work on weak supervision focuses on text or images. What about relational data? How can weak supervision be applied there?

Recipe for Creating Training Data

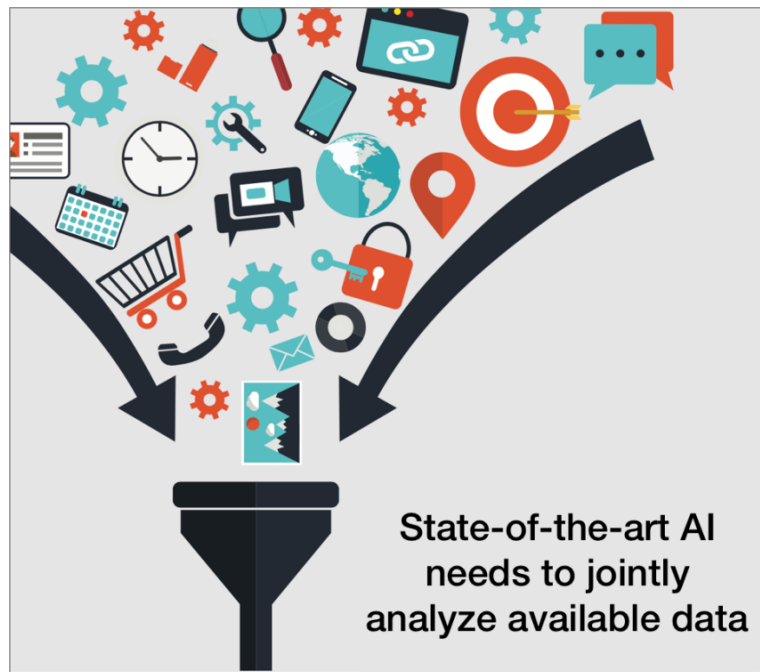
- Problem definition: Go beyond gold labels to noisy training data.
- Short answers
 - Transition from “gold” labels to “high-confidence” labels.
 - Modeling error rates is key. The notion of *data source* is different.
 - Need for debugging tools, bias detection, and recommendations of weak supervision signals.



Outline

- Part I. Introduction
- Part II. ML for DI
- Part III. DI for ML
 - Training data creation
 - Data cleaning
- Part IV. Conclusions and research directions

Successful ML requires Data Integration



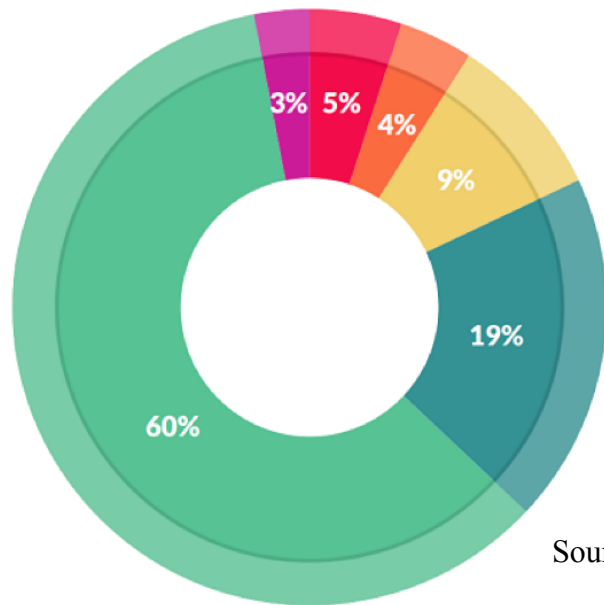
IMAGENET MovieLens



COCO is a large-scale object detection, segmentation, and captioning dataset.

Large collections of **manually curated** training data are necessary for progress in ML.

Noisy data is a bottleneck



What data scientists spend the most time doing

- Building training sets: 3%
- Cleaning and organizing data: 60%
- Collecting data sets; 19%
- Mining data for patterns: 9%
- Refining algorithms: 4%
- Other: 5%

Source: Crowdfunder

Cleaning and organizing data comprises 60% of the time spent on an analytics of AI project.

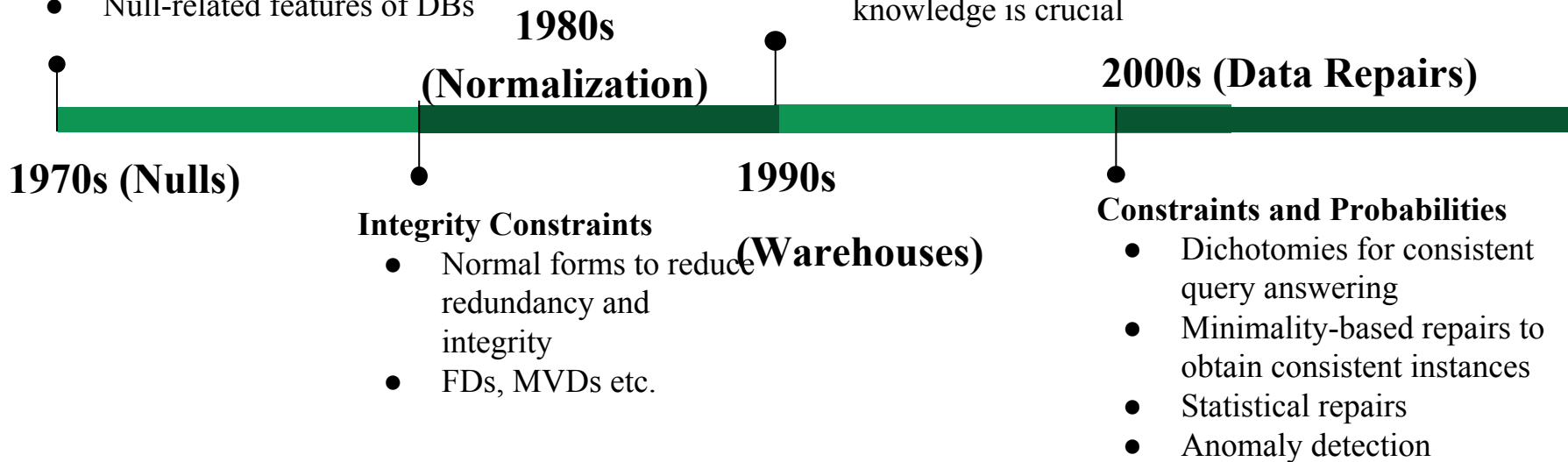
50 Years of Data Cleaning

E. F. Codd

- Understanding relations (installment #7).
FDT - Bulletin of ACM SIGMOD, 7(3):23–28, 1975.
- Null-related features of DBs

Data transforms

- Part of ETL
- Errors within a source and across sources
- Transformation workflows and mapping rules; domain-knowledge is crucial

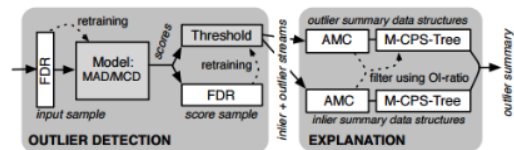
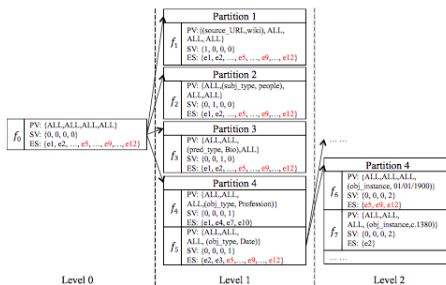
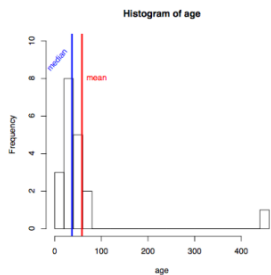


Where are we today?

Machine learning and statistical analysis are becoming more prevalent.

Error detection (Diagnosis)

- Anomaly detection [Chandola et al., ACM CSUR, 2009]
- Bayesian analysis (Data X-Ray) [Wang et al., SIGMOD'15]
- Outlier detection over streams (Macrobase) [Bailis et al., SIMGOD'17]

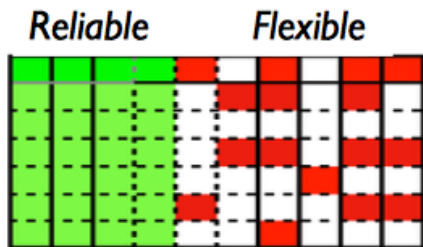


Where are we today?

Machine learning and statistical analysis are becoming more prevalent.

Data Repairing (Treatment)

- Classical ML (SCARE, ERACER) [Yakout et al., VLDB'11, SIGMOD'13, Mayfield et al., SIGMOD'10]
- Boosting [Krishan et al., 2017]
- Weakly-supervised ML (HoloClean) [Rekatsinas et al., VLDB'17]



Each cell is a random variable

Address	City	State	Zip
3465 S Morgan ST	Chicago	IL	60608
3465 S Morgan ST	Chicago	IL	60609
3465 S Morgan ST	Chicago	IL	60609
3465 S Morgan ST	Chicago	IL	60608

Constraints introduce correlations
c3: City, State, Address \rightarrow Zip

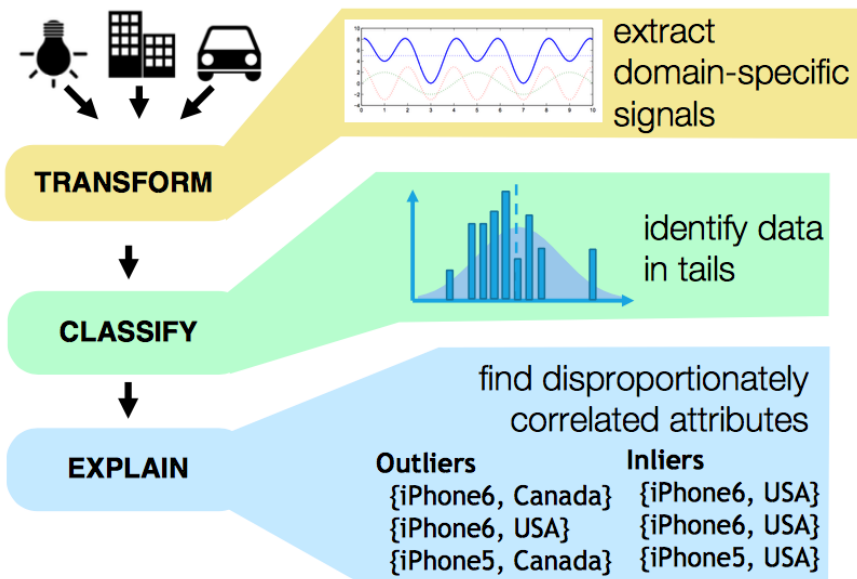
External data introduce evidence

Ext_Address	Ext_City	Ext_State	Ext_Zip
3465 S Morgan ST	Chicago	IL	60608





Error Detection: MacroBase [Bailis et al., SIGMOD'17]



[Figure by Kai Sheng Tai]

Streaming Feature Selection

Setup: Online learning of a classifier (e.g., LR)

Goal: Return top-k discriminative features

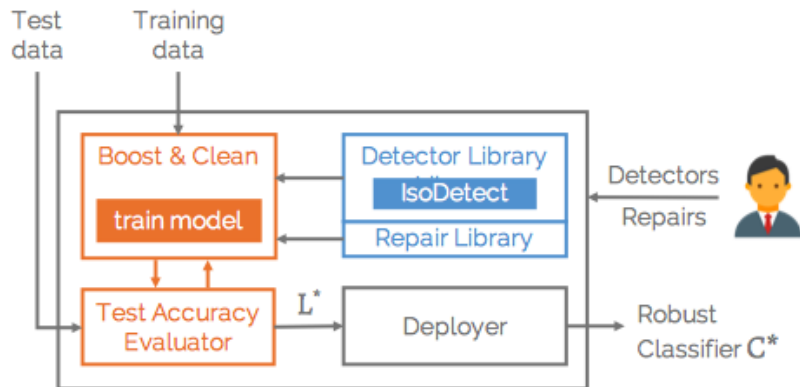
Weight-Median Sketch

Sketch of a classifier for fast updates and queries for estimates of each weight and comes with approximation guarantees

A data analytics tool that prioritizes attention in large datasets.

Code at: macrobase.stanford.edu

Data Repairing: BoostClean [Krishnan et al., 2017]



Ensemble learning for error detection and data repairing.

Relies on domain-specific detection and repairing.

Builds upon boosting to identify repairs that will maximize the performance improvement of a downstream classifier.

On-demand cleaning!

Scalable machine learning for data enrichment



 HoloClean

The logo for HoloClean features a stylized red 'H' composed of four teardrop-like shapes connected by a horizontal bar with two small circles. To the right of this icon is the word 'HoloClean' in a red, sans-serif font.

Code available at:

<http://www.holoclean.io>



Data Repairing: HoloClean [Rekatsinas et al., VLDB'17]

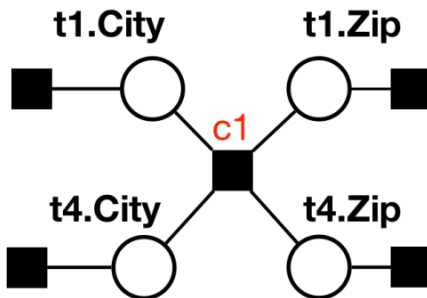
Each cell is a random variable

Value co-occurrences capture data statistics

Constraints introduce correlations

$c1: \text{Zip} \rightarrow \text{City}$

"Address= 3465 S Morgan St"



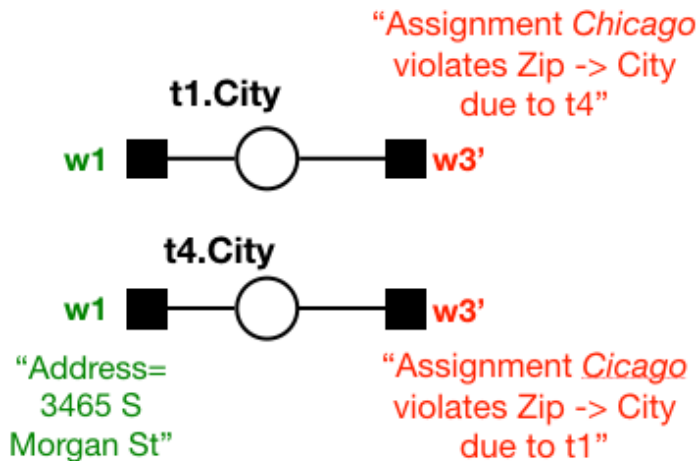
○ : Unknown (to be inferred) RV
■ : Factor (encodes correlations)

Holistic data cleaning framework: combines a variety of heterogeneous signals (e.g., integrity constraints, external knowledge, quantitative statistics)



Data Repairing: HoloClean [Rekatsinas et al., VLDB'17]

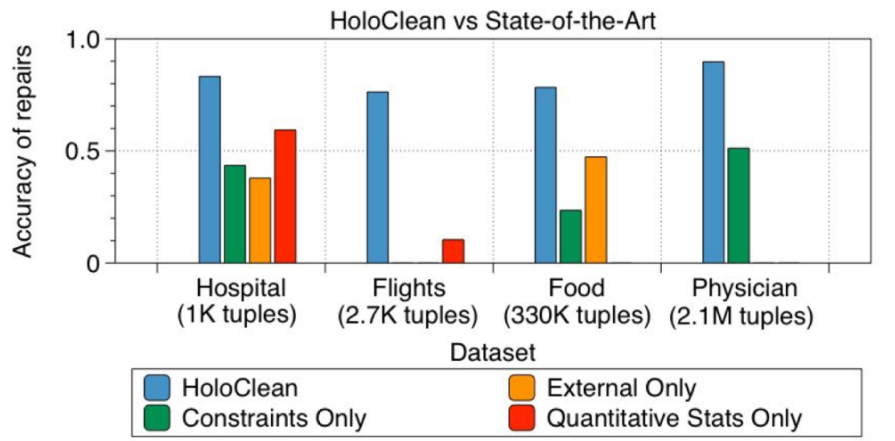
	Address	City	State	Zip
t1	3465 S Morgan ST	Chicago	IL	60608
t2	3465 S Morgan ST	Chicago	IL	60609
t3	3465 S Morgan ST	Chicago	IL	60609
t4	3465 S Morgan ST	Cicago	IL	60608



Scalable learning and inference: Hard constraints lead to complex and non-scalable models. Novel relaxation to features over individual cells.



Data Repairing: HoloClean [Rekatsinas et al., VLDB'17]



HoloClean is 2x more accurate. Competing methods either do not scale or perform no correct repairs.

HoloClean: our approach combining all signals and using inference

Holistic[Chu,2013]: state-of-the-art for constraints & minimality

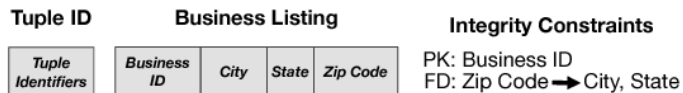
KATARA[Chu,2015]: state-of-the-art for external data

SCARE[Yakout,2013]: state-of-the-art ML & qualitative statistics

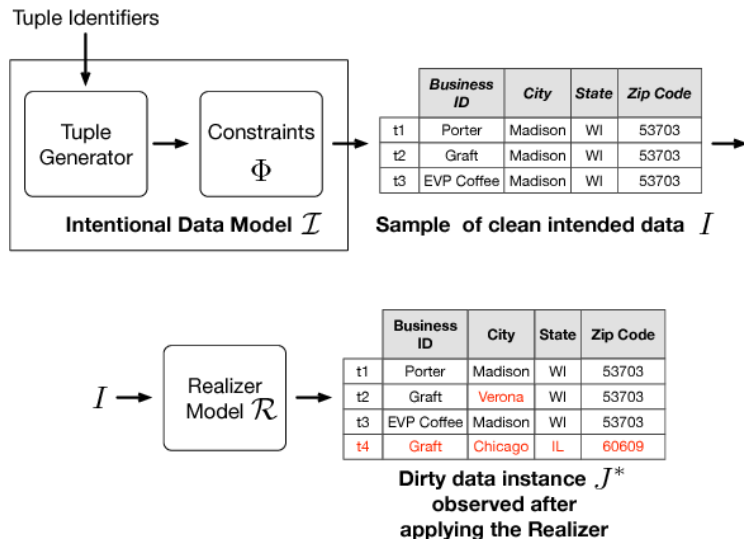
Probabilistic Unclean Databases [De Sa et al., 2018]

Unclean Database Generation

(A) Schema, Attribute Domain, and Constraint Specification



(B) The Two-Actor Generation Process



A two-actor noisy channel model for managing erroneous data.

Preprint: *A Formal Framework For Probabilistic Unclean Databases*

<https://arxiv.org/abs/1801.06750>

Challenges in Data Cleaning

- Error detection is still a challenge. To what extent is ML useful for error detection? Tuple-scoped approaches seem to be dominating. Is deep learning useful?
- We need a formal framework to describe when automated solutions are possible.
- A major bottleneck is the collection of training data. Can we leverage weak supervision and data augmentation more effectively?
- Limited end-to-end solutions. Data cleaning workloads (mixed relational and statistical workloads) pose unique scalability challenges.

Recipe for Data Cleaning

- Problem definition: **Detect and repair erroneous data.**

- Short answers

- **ML can help partly-automate cleaning. Domain expertise is still required.**
- **Scalability of ML-based data cleaning methods is a pressing challenge. Exciting systems research!**
- **We need more end-to-end systems!**

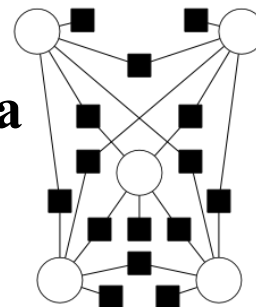
Each cell is a random variable

Address	City	State	Zip
3465 S Morgan ST	Chicago	IL	60608
3465 S Morgan ST	Chicago	IL	60609
3465 S Morgan ST	Chicago	IL	60609
3465 S Morgan ST	Chicago	IL	60609
3465 S Morgan ST	Chicago	IL	60608

Constraints introduce correlations
c3: City, State, Address \rightarrow Zip

External data introduce evidence

Ext_Address	Ext_City	Ext_State	Ext_Zip
3465 S Morgan ST	Chicago	IL	60608



Outline

- Part I. Introduction
- Part II. ML for DI
- Part III. DI for ML
 - Creating training data
 - Data cleaning
- **Part IV. Conclusions and research direction**

DI and ML: A Natural Synergy

- Data integration is one of the oldest problems in data management
- Transition from logic to probabilities revolutionized data integration
 - Probabilities allow us to reason about inherently noisy data
 - Similar to the AI-revolution in the 80s [<https://vimeo.com/48195434>]
- Modern machine learning and deep learning have the power to streamline DI

DI and ML: A Natural Synergy

- Data is bottleneck of modern ML and AI applications
- DI-related methods and algorithms have revolutionized the way supervision is performed.
 - Weak supervision signals are integrated into training datasets
- Data integration solutions (e.g., data cataloging solutions) can lead to cheaper collection of training data and more effective data enrichment

Opportunities for DI

One System vs. An Ecosystem: Every RBMS is a monolithic system. This paradigm has failed for DI. Tools for different DI tasks are prevalent. We need abstractions and execution frameworks for such ecosystems.

Humans-in-the-loop: DI tasks can be very complex. Is weak supervision the right approach to inject domain knowledge? What about quality evaluation?

Multi-modal DI: ML-based DI has focused on structured data with the exception of DI over images using crowdsourcing and some recent efforts that target textual data. DL is the de facto solution to reasoning about high dimensional data. Can is help develop unified DI solutions for visual, textual, and structured data?

Efficient Model Serving: This means efficient model serving. Many compute-intensive operations such as normalization and blocking are required. Featurization may also rely on compute-heavy tasks (e.g., computing string similarity). What is the role of pipelining and RDBMS-style optimizations?

Opportunities for ML

Data Catalogs: Data augmentation relies on data transformations performed on data records in a single dataset. How can we leverage data catalogs and data hubs to enable data augmentation go beyond a single dataset?

Valuable Data for ML applications: Our community has focused on assessing the value of data [Dong et al., VLDB'12, Koutris et al., JACM 2015]. These ideas are not pervasive to ML but if ML is to become a commodity [Jordan, 2018] we need methods to reason about the value of data.

DI for Benchmarks: Increasing efforts on creating manually curated benchmarks for ML. Current efforts rely on manual collection and curation. How can we leverage meta-data and existing DI solutions to automate such efforts?

“How reliable are our current measures of progress in machine learning?”
Do CIFAR-10 Classifiers Generalize to CIFAR-10?, Ben Recht et al., 2018



DI & ML as Synergy

- **ML for effective DI: AUTOMATION, AUTOMATION, AUTOMATION**
 - Automating DI tasks with training data
 - Ensemble learning and deep learning provide promising solutions
 - Better understanding of semantics by neural network
- **DI for effective ML: DATA, DATA, DATA**
 - The software 2.0 stack is data hungry
 - Create large-scale training datasets from different sources
 - Cleaning of data used for training

Thank you!

References Part I: Introduction

- Bengio, Y., Goodfellow, I.J. & Courville, A., 2015. Deep learning. *Nature*, 521(7553), pp.436–444.
- Bishop, C.M., 2016. *Pattern Recognition and Machine Learning*, Springer New York.
- Doan, A., Halevy, A.Y. & Ives, Z.G., 2012. *Principles of Data Integration*, Morgan Kaufmann.
- Domingos, P., 2012. A Few Useful Things to Know About Machine Learning. *Communications of the ACM*, 55(10), pp.78–87.
- Dong, X. et al., 2014. Knowledge Vault: A Web-scale Approach to Probabilistic Knowledge Fusion. In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. KDD '14. New York, NY, USA: ACM, pp. 601–610.
- Dong, X.L. & Srivastava, D., 2015. Big data integration. *Synthesis Lectures on Data Management*, 7(1), pp.1–198.
- Dong, X.L. & Srivastava, D., 2013. Big Data Integration. *Proceedings of the VLDB Endowment International Conference on Very Large Data Bases*, 6(11), pp.1188–1189.
- Getoor, L. & Machanavajjhala, A., 2012. Entity resolution: theory, practice & open challenges. *PVLDB*, 5(12), pp.2018–2019.
- Goodfellow, I. et al., 2016. *Deep learning*, MIT press Cambridge.
- Halevy, A., Norvig, P. & Pereira, F., 2009. The Unreasonable Effectiveness of Data. *IEEE intelligent systems*, 24(2), pp.8–12.
- Konda, P. et al., 2016. Magellan: Toward Building Entity Matching Management Systems. *PVLDB*, 9(12), pp.1197–1208.

References Part I: Introduction

- Kumar, A., Boehm, M. & Yang, J., 2017. Data Management in Machine Learning: Challenges, Techniques, and Systems. In *Proceedings of the 2017 ACM International Conference on Management of Data*. SIGMOD '17. New York, NY, USA: ACM, pp. 1717–1722.
- Lockard, C. et al., 2018. CERES: Distantly Supervised Relation Extraction from the Semi-Structured Web. *arXiv [cs.AI]*. Available at: <http://arxiv.org/abs/1804.04635>.
- Mohri, M., Rostamizadeh, A. & Talwalkar, A., 2012. *Foundations of Machine Learning*, MIT Press.
- Polyzotis, N. et al., 2017. Data Management Challenges in Production Machine Learning. In *Proceedings of the 2017 ACM International Conference on Management of Data*. SIGMOD '17. New York, NY, USA: ACM, pp. 1723–1726.
- Ratner, A. et al., 2017. Snorkel: Rapid Training Data Creation with Weak Supervision. *PVLDB*, 11(3), pp.269–282.
- Rekatsinas, T. et al., 2017. HoloClean: Holistic Data Repairs with Probabilistic Inference. *PVLDB*, 10(11), pp.1190–1201.
- Wu, S. et al., 2018. Fonduer: Knowledge Base Construction from Richly Formatted Data. In *Proceedings of the 2018 International Conference on Management of Data*. ACM, pp. 1301–1316.
- Zheng, G. et al., 2018. OpenTag: Open Attribute Value Extraction from Product Profiles. In *KDD*. Available at: <https://people.mpi-inf.mpg.de/~smukherjee/research/OpenTag-KDD18.pdf>.

References Part II: Entity Linkage

- Bhattacharya, I. & Getoor, L., 2006. A latent dirichlet model for unsupervised entity resolution. In *SDM*. SIAM, pp. 47–58.
- Das, S. et al., 2017. Falcon: Scaling Up Hands-Off Crowdsourced Entity Matching to Build Cloud Services. In *Sigmod*. pp. 1431–1446.
- Doan, A. et al., 2017. Human-in-the-Loop Challenges for Entity Matching: A Midterm Report. In *Proceedings of the 2nd Workshop on Human-In-the-Loop Data Analytics, HILDA@SIGMOD 2017, Chicago, IL, USA, May 14, 2017*. pp. 12:1–12:6.
- Fellegi, I.P. & Sunter, A.B., 1969. A Theory for Record Linkage. *Journal of the Americal Statistical Association*, 64(328), pp.1183–1210.
- Getoor, L. & Machanavajjhala, A., 2012. Entity resolution: theory, practice & open challenges. *PVLDB*, 5(12), pp.2018–2019.
- Gokhale, C. et al., 2014. Corleone: Hands-off Crowdsourcing for Entity Matching. In *Proceedings of the 2014 ACM SIGMOD International Conference on Management of Data*. SIGMOD '14. New York, NY, USA: ACM, pp. 601–612.
- Hassanzadeh, O. et al., 2009. Framework for Evaluating Clustering Algorithms in Duplicate Detection. *PVLDB*, 2(1), pp.1282–1293.
- Ji, H., 2014. Entity Linking and Wikification Reading List. Available at: <http://nlp.cs.rpi.edu/kbp/2014/elreading.html>.
- Konda, P. et al., 2016. Magellan: Toward Building Entity Matching Management Systems. *PVLDB*, 9(12), pp.1197–1208.
- Kopcke, H., Thor, A. & Rahm, E., 2010. Evaluation of entity resolution approaches on real-world match problems. *PVLDB*, 3(1), pp.484–493.

References Part II: Entity Linkage

- Mudgal, S. et al., 2018. Deep Learning for Entity Matching: A Design Space Exploration. In *Proceedings of the 2018 International Conference on Management of Data*. ACM, pp. 19–34.
- Pujara, J. & Getoor, L., 2016. Generic Statistical Relational Entity Resolution in Knowledge Graphs. In *AAAI*.
- Rakshit Trivedi, Bunyamin Sisman, Xin Luna Dong, Christos Faloutsos, Jun Ma and Hongyuan Zha., LinkNBed: Multi-Graph Representation Learning with Entity Linkage. In *56th Annual Meeting of the Association for Computational Linguistics*. ACL.
- Sarawagi, S. & Bhamidipaty, A., 2002. Interactive deduplication using active learning. In *SIGKDD*.
- Singla, P. & Domingos, P., 2006. Entity Resolution with Markov Logic. In *ICDM*. Washington, DC, USA: IEEE Computer Society, pp. 572–582.
- Stonebraker, M. et al., 2013. Data Curation at Scale: The Data Tamer System. In *CIDR*.
- Verroios, V., Garcia-Molina, H. & Papakonstantinou, Y., 2017. Waldo: An Adaptive Human Interface for Crowd Entity Resolution. In *Proceedings of the 2017 ACM International Conference on Management of Data, SIGMOD Conference 2017, Chicago, IL, USA, May 14-19, 2017*. pp. 1133–1148.

References Part II: Data Extraction

- Das, R. et al., 2017. Chains of reasoning over entities, relations, and text using recurrent neural networks. In *EACL*.
- Dong, X. et al., 2014. Knowledge Vault: A Web-scale Approach to Probabilistic Knowledge Fusion. In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. KDD '14. New York, NY, USA: ACM, pp. 601–610.
- Dong, X.L., 2017. Challenges and Innovations in Building a Product Knowledge Graph. In *AKBC*.
- Gulhane, P. et al., 2011. Web-scale information extraction with vertex. In *2011 IEEE 27th International Conference on Data Engineering*. pp. 1209–1220.
- He, R. et al., 2017. An Unsupervised Neural Attention Model for Aspect Extraction. In *ACL*.
- Hoffmann, R. et al., 2011. Knowledge-based Weak Supervision for Information Extraction of Overlapping Relations. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies - Volume 1*. HLT '11. Stroudsburg, PA, USA: Association for Computational Linguistics, pp. 541–550.
- Limaye, G., Sarawagi, S. & Chakrabarti, S., 2010. Annotating and Searching Web Tables Using Entities, Types and Relationships. *Proceedings of the VLDB Endowment International Conference on Very Large Data Bases*, 3(1-2), pp.1338–1347.
- Lockard, C. et al., 2018. CERES: Distantly Supervised Relation Extraction from the Semi-Structured Web. *arXiv [cs.AI]*. Available at: <http://arxiv.org/abs/1804.04635>.

References Part II: Data Extraction

- Mintz, M. et al., 2009. Distant supervision for relation extraction without labeled data. In *ACL*.
- Mitchell, T. et al., 2018. Never-ending Learning. *Communications of the ACM*, 61(5), pp.103–115.
- Neelakantan, A., Roth, B. & McCallum, A., 2015. Compositional vector space models for knowledge base completion. In *ACL*.
- Riedel, S. et al., 2013. Relation Extraction with Matrix Factorization and Universal Schemas. In *HLT-NAACL*.
- Shin, J. et al., 2015. Incremental Knowledge Base Construction Using DeepDive. *Proceedings of the VLDB Endowment International Conference on Very Large Data Bases*, 8(11), pp.1310–1321.
- Wu, S. et al., 2018. Fondue: Knowledge Base Construction from Richly Formatted Data. In *Proceedings of the 2018 International Conference on Management of Data*. ACM, pp. 1301–1316.
- Zhang, C. et al., 2017. DeepDive: Declarative Knowledge Base Construction. *CACM*, 60(5), pp.93–102.

References Part II: Data Fusion

- Dawid, A.P. & Skene, A.M., 1979. Maximum Likelihood Estimation of Observer Error-Rates Using the EM Algorithm. *Journal of the Royal Statistical Society. Series C, Applied statistics*, 28(1), pp.20–28.
- Dong, X.L. et al., 2014. From Data Fusion to Knowledge Fusion. *PVLDB*.
- Dong, X.L. et al., 2015. Knowledge-based Trust: Estimating the Trustworthiness of Web Sources. *Proceedings of the VLDB Endowment International Conference on Very Large Data Bases*, 8(9), pp.938–949.
- Dong, X.L. & Naumann, F., 2009. Data Fusion: Resolving Data Conflicts for Integration. *Proceedings of the VLDB Endowment International Conference on Very Large Data Bases*, 2(2), pp.1654–1655.
- Gao, J. et al., 2016. Mining Reliable Information from Passively and Actively Crowdsourced Data. In *KDD*. pp. 2121–2122.
- Jaffe, A., Nadler, B. & Kluger, Y., 2015. Estimating the accuracies of multiple classifiers without labeled data. In *Artificial Intelligence and Statistics*. Artificial Intelligence and Statistics. pp. 407–415.
- Li, H., Yu, B. & Zhou, D., 2013. Error rate analysis of labeling by crowdsourcing. In *ICML Workshop: Machine Learning Meets Crowdsourcing*. Atlanta, Georgia, USA.
- Li, Q. et al., 2014. A Confidence-aware Approach for Truth Discovery on Long-tail Data. *Proceedings of the VLDB Endowment International Conference on Very Large Data Bases*, 8(4), pp.425–436.

References Part II: Data Fusion

- Li, X. et al., 2013. Truth Finding on the Deep Web: Is the Problem Solved? *PVLDB*, 6(2).
- Li, Y. et al., 2016. A Survey on Truth Discovery. *SIGKDD Explor. Newsl.*, 17(2), pp.1–16.
- Nickel, M. et al., 2016. A Review of Relational Machine Learning for Knowledge Graphs. *Proceedings of the IEEE*, 104(1), pp.11–33.
- Pasternack, J. & Roth, D., 2010. Knowing what to believe (when you already know something). In *COLING*. pp. 877–885.
- Platanios, E. A., Dubey, A., & Mitchell, T. (2016, June). Estimating accuracy from unlabeled data: A bayesian approach. In *International Conference on Machine Learning*(pp. 1416-1425).
- Rekatsinas, T. et al., 2017. SLiMFast: Guaranteed Results for Data Fusion and Source Reliability. In *Proceedings of the 2017 ACM International Conference on Management of Data*. SIGMOD '17. New York, NY, USA: ACM, pp. 1399–1414.
- Shaham, U. et al., 2016. A Deep Learning Approach to Unsupervised Ensemble Learning. In *International Conference on Machine Learning*. International Conference on Machine Learning. pp. 30–39.
- Wang, Q. et al., 2017. Knowledge Graph Embedding: A Survey of Approaches and Applications. *IEEE transactions on knowledge and data engineering*, 29(12), pp.2724–2743.
- Yin, X., Han, J. & Yu, P.S., 2007. Truth discovery with multiple conflicting information providers on the web. In *Proceedings of the 13th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, pp. 1048–1052.

References Part II: Data Fusion

Zhang, Y. et al., 2014. Spectral Methods meet EM: A Provably Optimal Algorithm for Crowdsourcing. In Z. Ghahramani et al., eds.

Advances in Neural Information Processing Systems 27. Curran Associates, Inc., pp. 1260–1268.

Zhao, B. et al., 2012. A Bayesian Approach to Discovering Truth from Conflicting Sources for Data Integration. *Proceedings of the VLDB*

Endowment International Conference on Very Large Data Bases, 5(6), pp.550–561.

References Part III: Training Data Creation

- Chapelle, O., Scholkopf, B. & Eds., A.Z., 2009. Semi-Supervised Learning (Chapelle, O. et al., Eds.; 2006) [Book reviews]. *IEEE transactions on neural networks / a publication of the IEEE Neural Networks Council*, 20(3), pp.542–542.
- Dawid, A.P. & Skene, A.M., 1979. Maximum Likelihood Estimation of Observer Error-Rates Using the EM Algorithm. *Journal of the Royal Statistical Society. Series C, Applied statistics*, 28(1), pp.20–28.
- Mintz, M. et al., 2009. Distant supervision for relation extraction without labeled data. In *ACL*.
- Mitchell, T., 2017. Learning from Limited Labeled Data (But a Lot of Unlabeled Data). Available at: https://lld-workshop.github.io/slides/tom_mitchell_lld.pdf.
- Platanios, E.A., Dubey, A. & Mitchell, T., 2016. Estimating Accuracy from Unlabeled Data: A Bayesian Approach. In *International Conference on Machine Learning*. International Conference on Machine Learning. pp. 1416–1425.
- Ratner, A. et al., 2017. Snorkel: Rapid Training Data Creation with Weak Supervision. *PVLDB*, 11(3), pp.269–282.
- Ratner, A.J. et al., 2016. Data programming: Creating large training sets, quickly. In *Advances in Neural Information Processing Systems*. pp. 3567–3575.
- Raykar, V.C. et al., 2010. Learning From Crowds. *Journal of machine learning research: JMLR*, 11, pp.1297–1322.
- Recht, B. et al., 2018. Do CIFAR-10 Classifiers Generalize to CIFAR-10? *arXiv [cs.LG]*. Available at: <http://arxiv.org/abs/1806.00451>.

References Part III: Training Data Creation

- Roth, B. & Klakow, D., 2013. Combining generative and discriminative model scores for distant supervision. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*. pp. 24–29.
- Russell, S. & Stefano, E., 2017. Label-free supervision of neural networks with physics and domain knowledge. *Proceedings of AAAI*.
- Salimans, T. et al., 2016. Improved techniques for training gans. In *Advances in Neural Information Processing Systems*. pp. 2234–2242.
- Schapire, R.E. & Freund, Y., 2012. Boosting: Foundations and Algorithms. Adaptive computation and machine learning.

References Part III: Data Cleaning

- Bailis, P. et al., 2017. MacroBase: Prioritizing Attention in Fast Data. In *Proceedings of the 2017 ACM International Conference on Management of Data*. SIGMOD '17. New York, NY, USA: ACM, pp. 541–556.
- Chu, X. et al., 2016. Data Cleaning: Overview and Emerging Challenges. In *Proceedings of the 2016 International Conference on Management of Data*. SIGMOD '16. New York, NY, USA: ACM, pp. 2201–2206.
- Chandola, V., Banerjee, A. & Kumar, V., 2009. Anomaly Detection: A Survey. *ACM Comput. Surv.*, 41(3), pp.15:1–15:58.
- Galhardas, H. et al., 2001. Declarative data cleaning: Language, model, and algorithms. In *VLDB*. pp. 371–380.
- Hellerstein, J.M., 2008. Quantitative data cleaning for large databases. *Statistical journal of the United Nations Economic Commission for Europe*. Available at: <http://db.cs.berkeley.edu/jmh/papers/cleaning-unece.pdf>.
- Ilyas, I.F., 2016. Effective Data Cleaning with Continuous Evaluation. *IEEE Data Eng. Bull.*, 39, pp.38–46.
- Krishnan, S. et al., 2016. ActiveClean: Interactive Data Cleaning for Statistical Modeling. *Proceedings of the VLDB Endowment International Conference on Very Large Data Bases*, 9(12), pp.948–959.
- Krishnan, S. et al., 2017. BoostClean: Automated Error Detection and Repair for Machine Learning. *arXiv [cs.DB]*. Available at: <http://arxiv.org/abs/1711.01299>.

References Part III: Data Cleaning

- Mayfield, C., Neville, J. & Prabhakar, S., 2010. ERACER: A Database Approach for Statistical Inference and Data Cleaning. In *Proceedings of the 2010 ACM SIGMOD International Conference on Management of Data*. SIGMOD '10. New York, NY, USA: ACM, pp. 75–86.
- Rekatsinas, T. et al., 2017. HoloClean: Holistic Data Repairs with Probabilistic Inference. *PVLDB*, 10(11), pp.1190–1201.
- Wang, X., Dong, X.L. & Meliou, A., 2015. Data X-Ray: A Diagnostic Tool for Data Errors. In *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data*. SIGMOD '15. New York, NY, USA: ACM, pp. 1231–1245.
- Yakout, M., Berti-Équille, L. & Elmagarmid, A.K., 2013. Don'T Be SCARED: Use SCalable Automatic REpairing with Maximal Likelihood and Bounded Changes. In *Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data*. SIGMOD '13. New York, NY, USA: ACM, pp. 553–564.