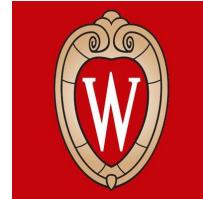
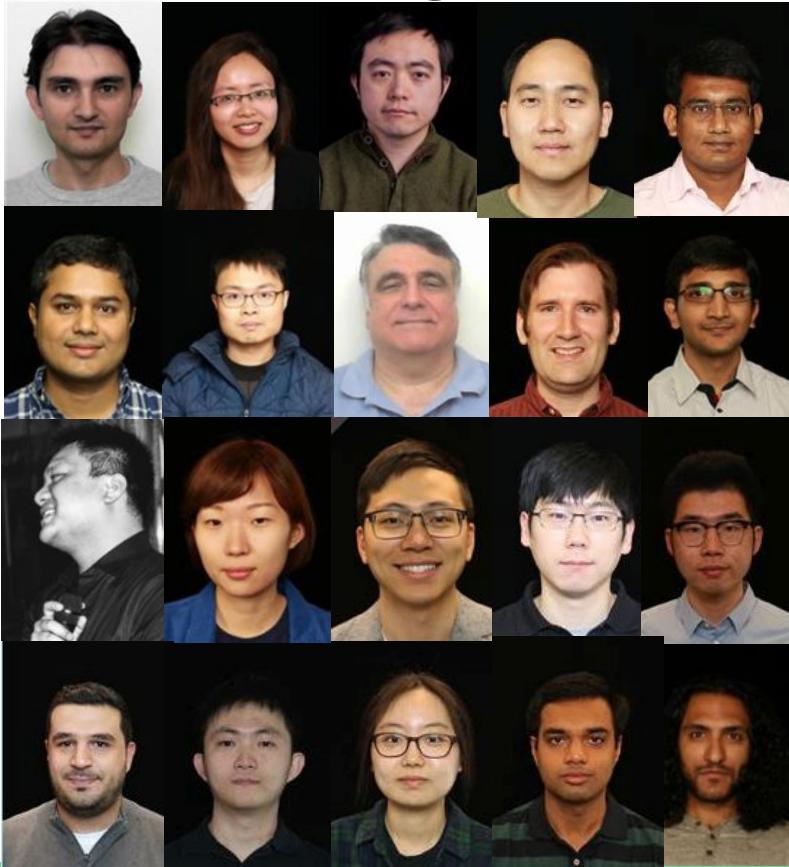


Data Integration and Machine Learning: A Natural Synergy

Xin Luna Dong @ Amazon.com
Theo Rekatsinas @ UW-Madison
<http://dataintegration.ai>

Acknowledgement

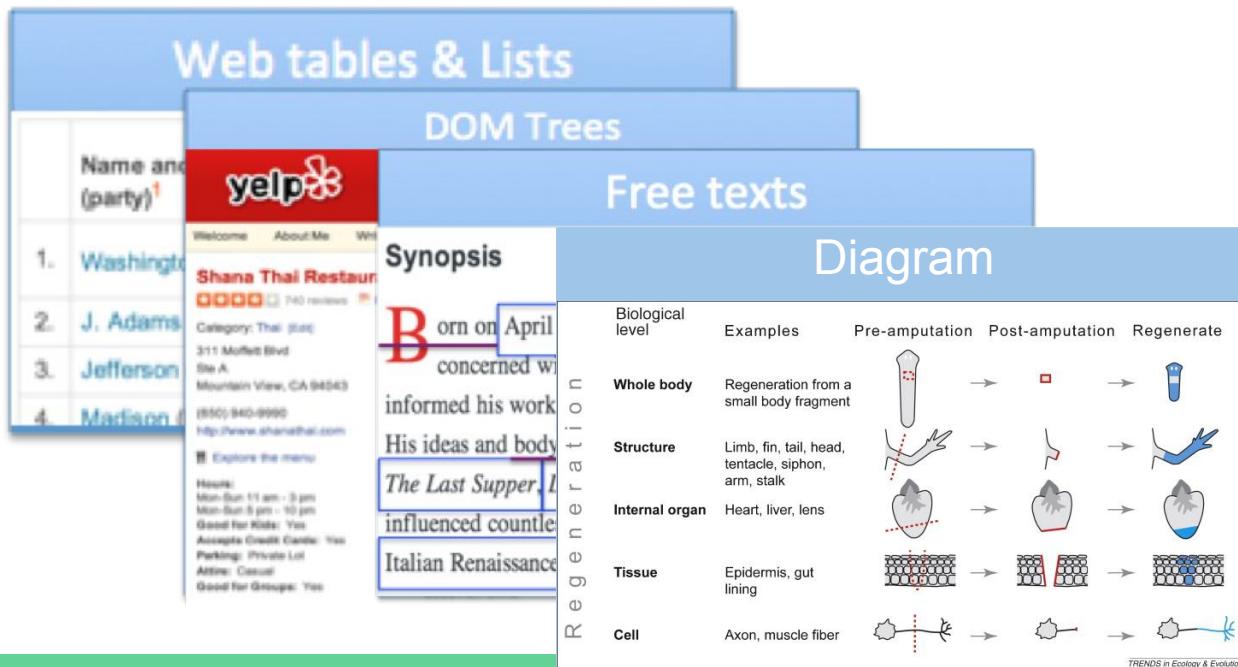


What is Data Integration?

- **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



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

SEE RANK

Anahí 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

Data Extraction

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

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

Actress | Music Department | Soundtrack

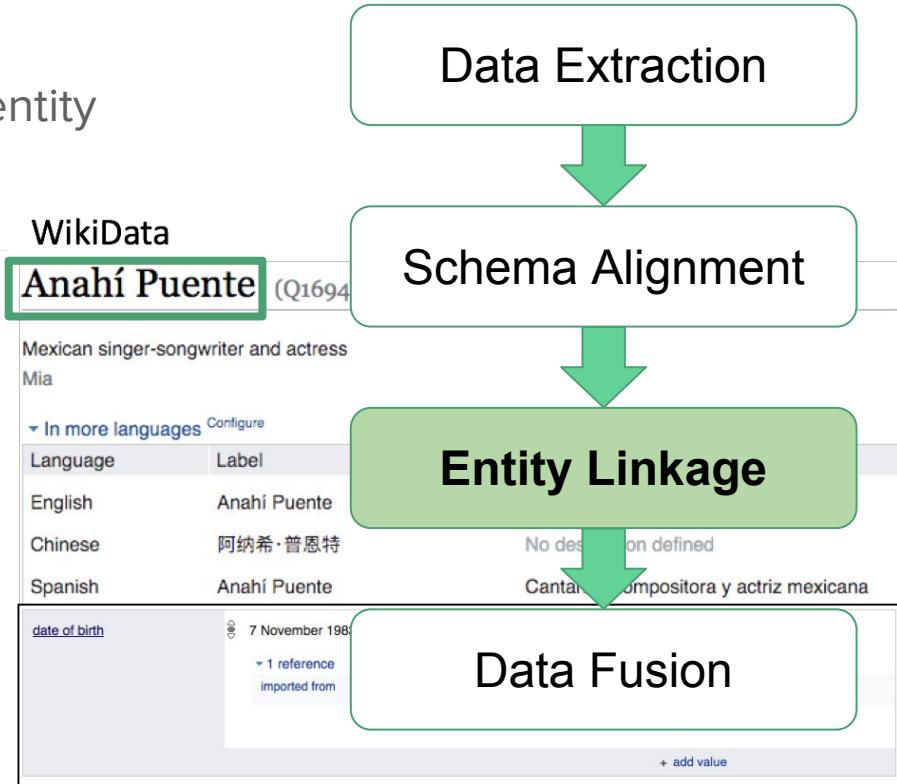
SEE RANK

Anahí 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. Anahí 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](#)



Why is Data Integration Hard?

- Heterogeneity everywhere
 - Conflicting values

IMDB



Anahí
Actress | Music Department | Soundtrack

SEE RANK

Anahí 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. Anahí 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

Data Extraction

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

Schema Alignment

Entity Linkage

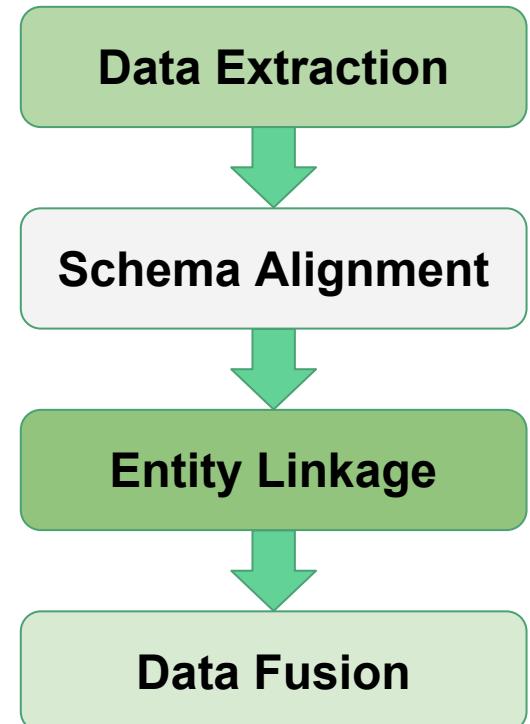
No description defined
Cantante, compositora y actriz mexicana

Data Fusion

+ add value

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



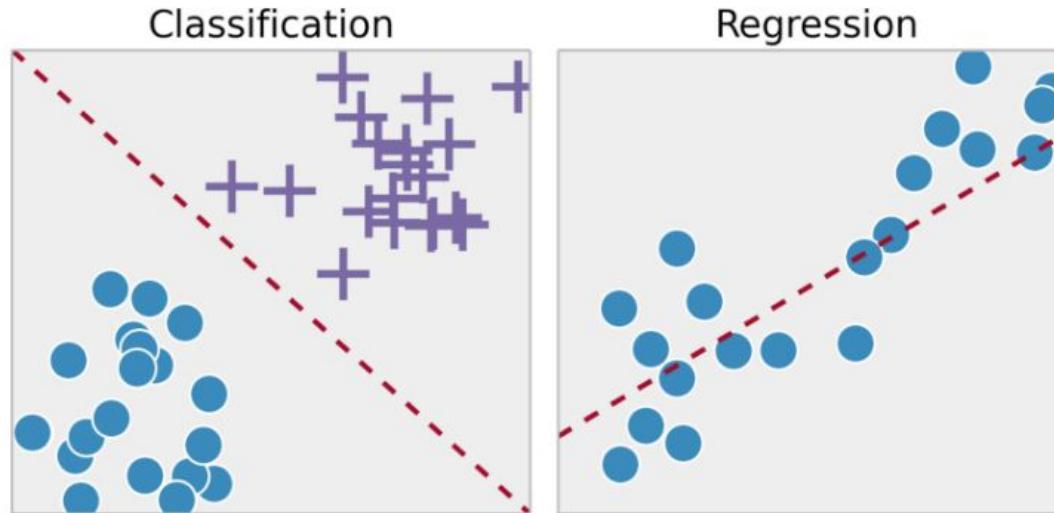
What is Machine Learning?

- **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

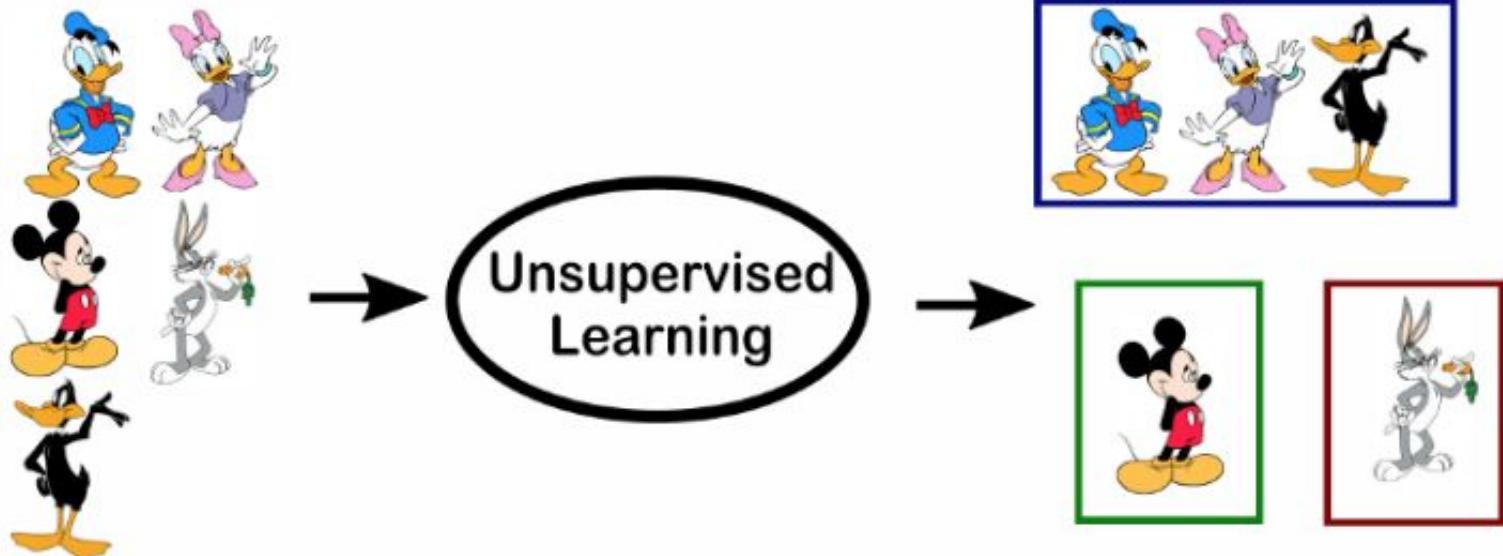
Two Main Types of Machine Learning

- Supervised learning: learn by examples



Two Main Types of Machine Learning

- Unsupervised learning: find structure w/o examples

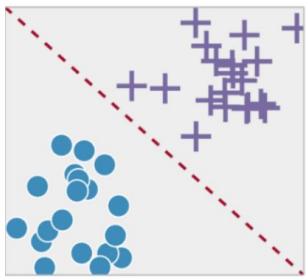
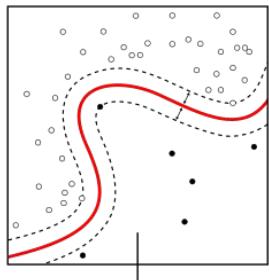
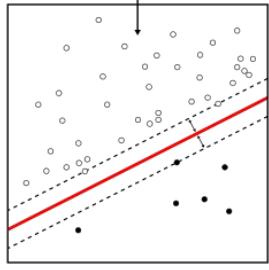
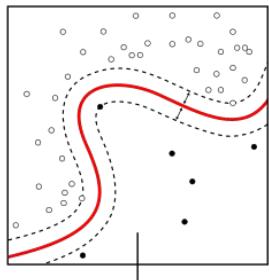
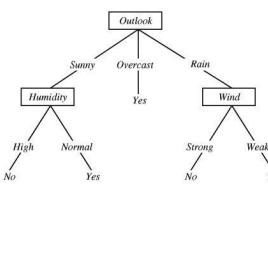
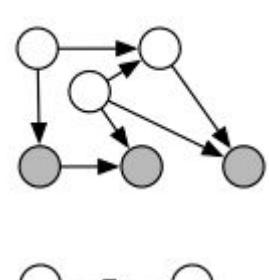
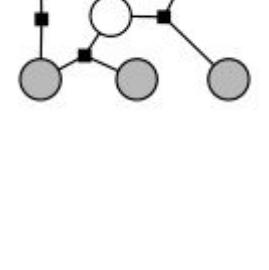
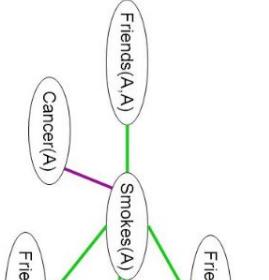
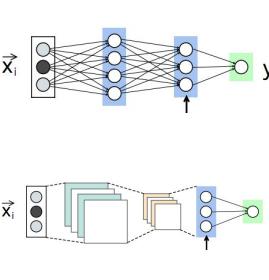
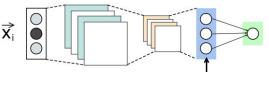
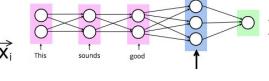


Two Main Types of Machine Learning

- 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

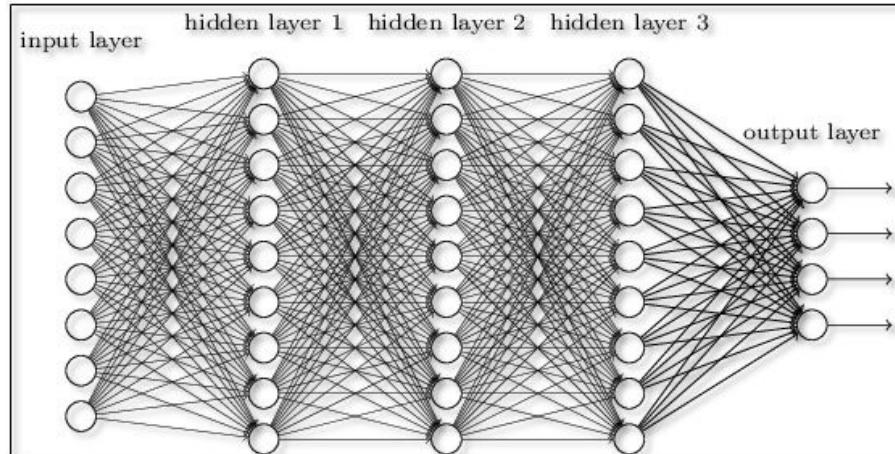
Techniques for Supervised ML

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, CNN, RNN
  			 		  

Colin's 7 minute intro to bidirectional LSTMs

Everyone loves neural networks

- They have the word “neural” in them
- They have the word “network” in them
- You can draw fancy diagrams
-



Everyone loves neural networks

- They have the word “neural” in them
- They have the word “network” in them
- You can draw fancy diagrams
- This thing:



Everyone loves neural networks

- They have the word “neural” in them
- They have the word “network” in them
- You can draw fancy diagrams
- This thing:



- They are good at modeling high-dimensional non-linear functions and building sophisticated representations of complex sensory data

What's a neural network?

x
↑
Input vector

What's a neural network?

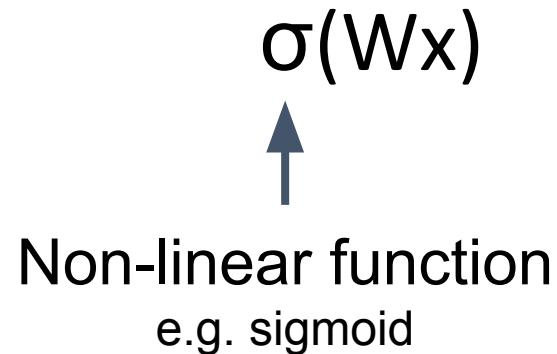
Wx



Weight matrix

What's a neural network?

$\sigma(Wx)$



↑

Non-linear function
e.g. sigmoid

What's a neural network?

$$h = \sigma(Wx)$$



Hidden layer unit

What's a neural network?

$$h = \sigma(Wx)$$

$$y = f(h)$$



Differentiable function
e.g. softmax

Let's use neural networks for language!

Classify sentences!

Tag parts of speech!

Find entity names!

Extract relations!

Let's use neural networks for language!

PROBLEM:

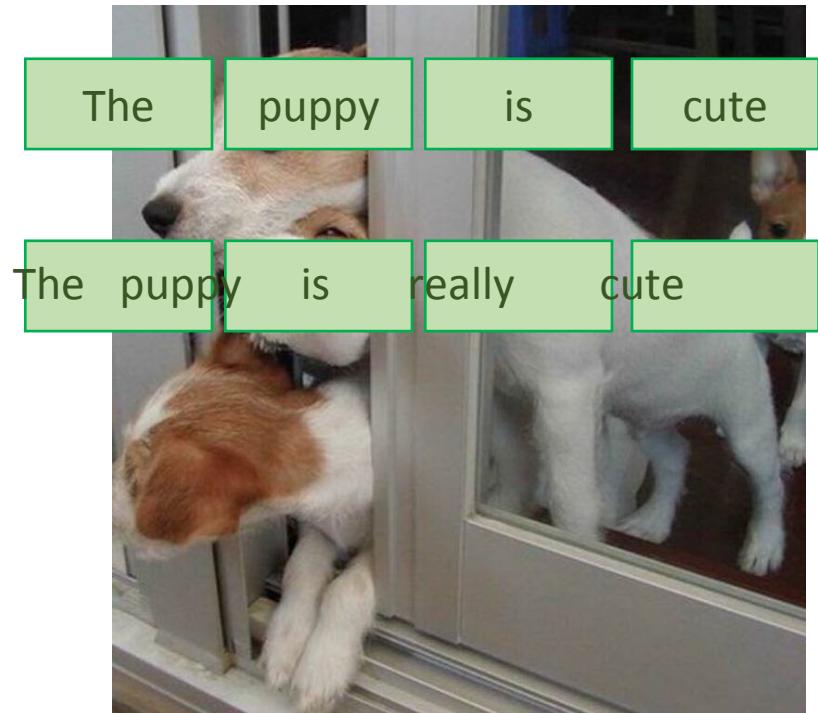
- Fixed # of weights
- Fixed # of features
- Fixed size of input



The puppy is cute.

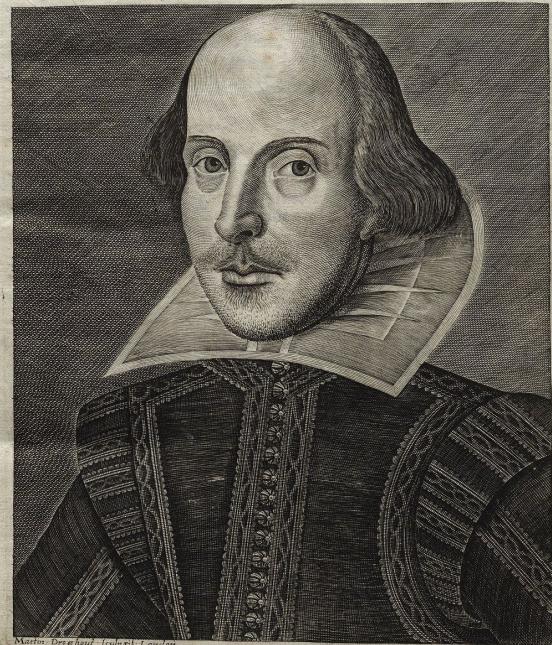
The puppy is really cute.

Seriously, did you see this puppy, look how cute it is!



MR. WILLIAM
SHAKESPEARES
COMEDIES,
HISTORIES, &
TRAGEDIES.

Published according to the True Originall Copies.



Martin Droeshout sculpsit London.

L O N D O N
Printed by Isaac Iaggard, and Ed. Blount. 1623.

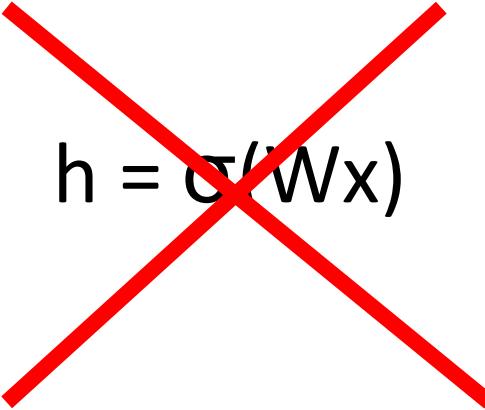


- Sentence classification
 - Sentences have variable length!
- Word tagging
 - Fixed-width window around word misses full context

What's a neural network?

$$h = \sigma(Wx)$$
A large, solid red 'X' is drawn across the center of the slide, intersecting the text 'h = σ(Wx)'.

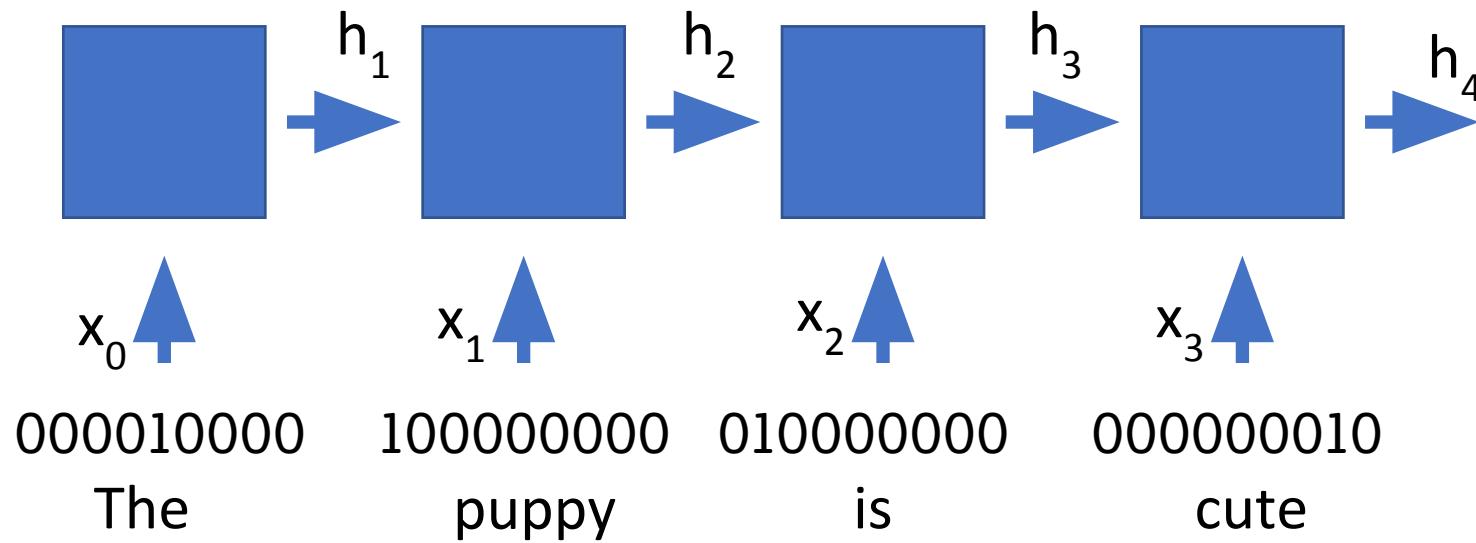
What's a recurrent neural network?

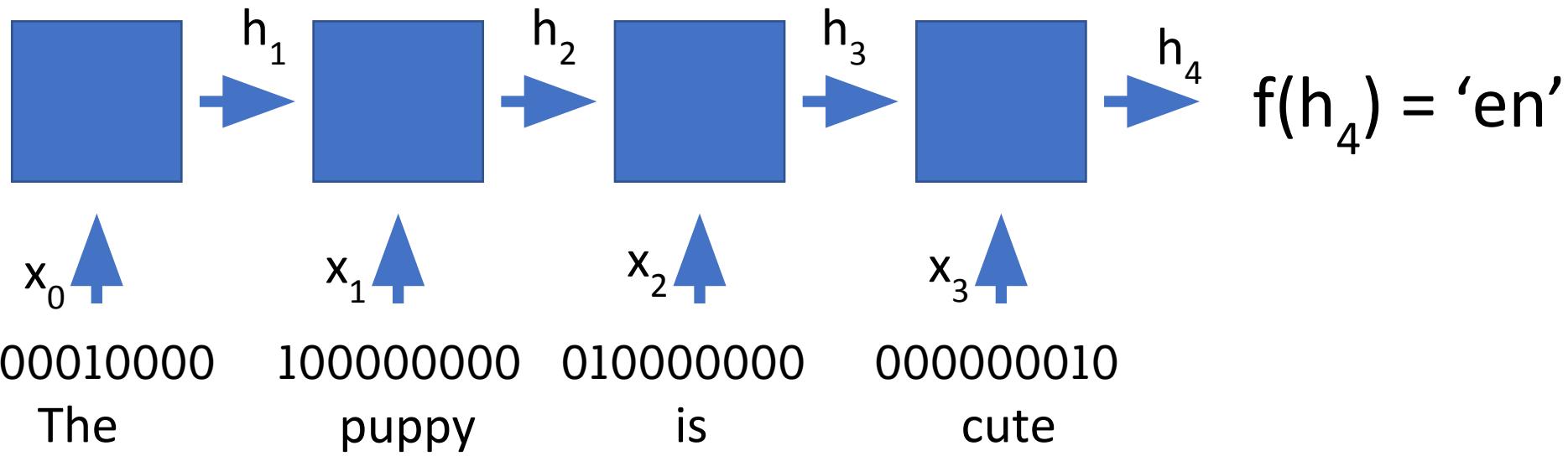
$$h = \sigma(Wx)$$


What's a recurrent neural network?

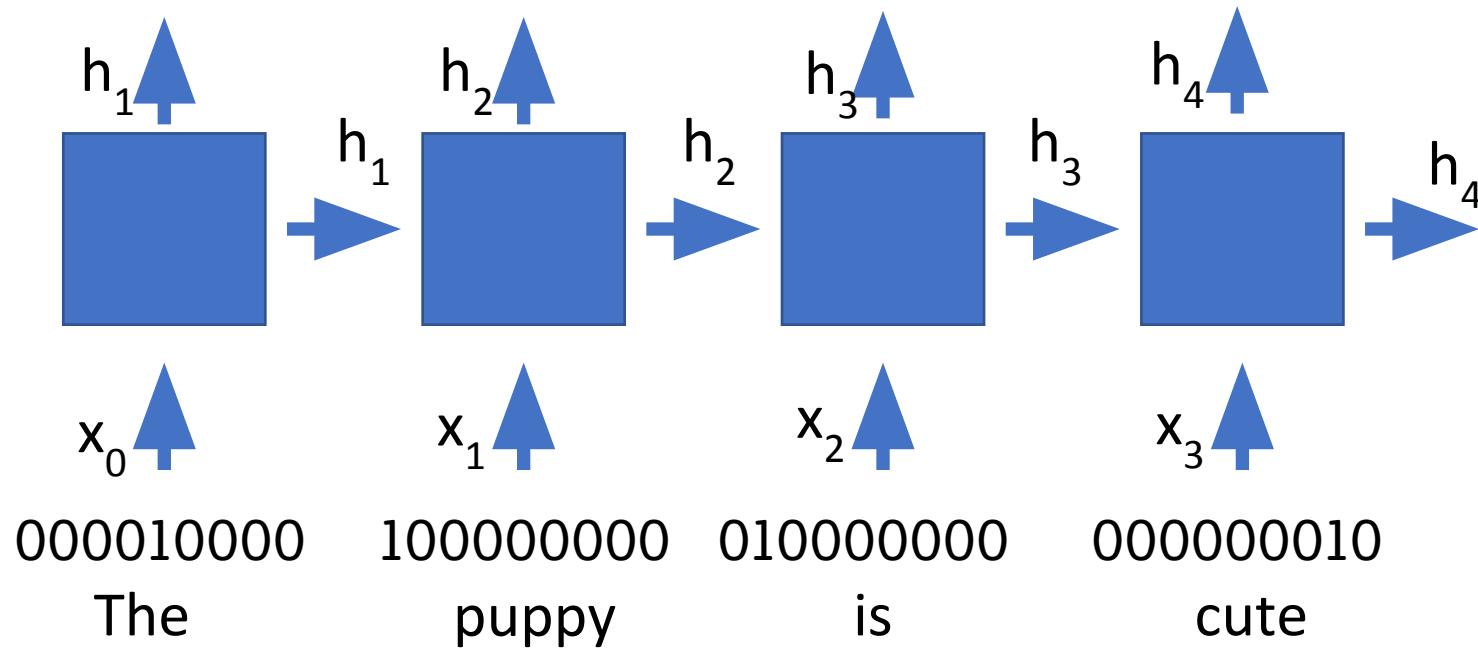
$$h = \sigma(Wx)$$

$$h_t = \sigma(Wx_t + Uh_{t-1})$$

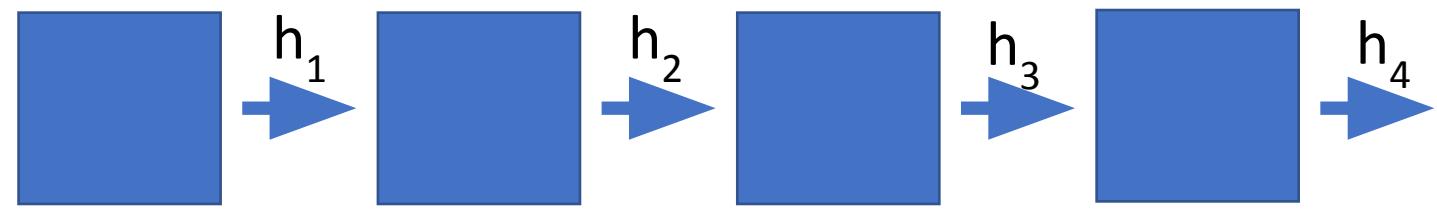




$f(h_1) \rightarrow$ def. art. $f(h_2) \rightarrow$ noun $f(h_3) \rightarrow$ verb $f(h_4) \rightarrow$ adjective



BIDIRECTIONAL



x_0

000010000

The

x_1

100000000

puppy

x_2

010000000

is

x_3

000000010

cute



x_3

000010000

The

x_2

100000000

puppy

x_1

010000000

is

x_0

000000010

cute

Bidirectional RNN

Concatenate the two hidden representations to produce the bidirectional representation

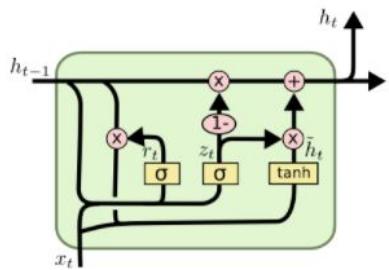
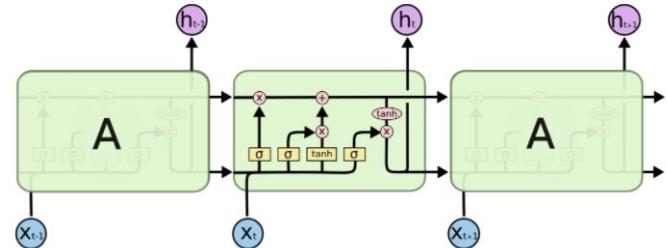
Full sentence: $f([h_4, h'_4])$

Individual word: $f([h_i, h'_{n-i}])$

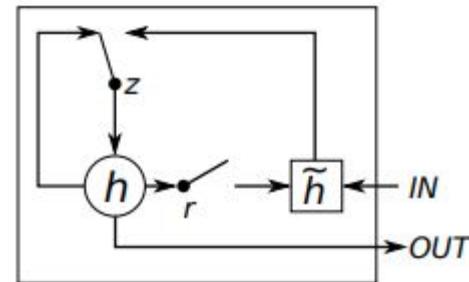
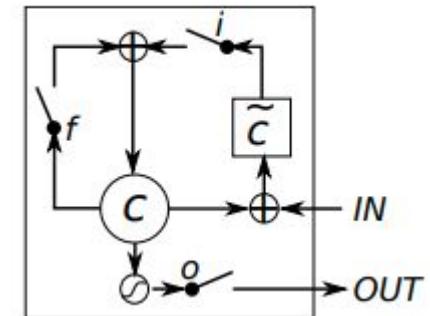
Okay, so RNNs are great!

- Except they don't work
 - (at least for long-term dependencies)
- Vanishing / exploding gradient

LSTM (Long Short Term Memory)

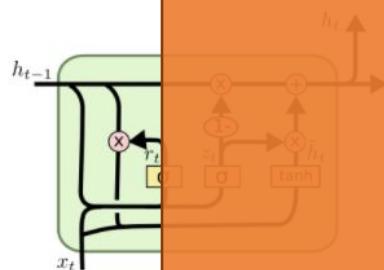
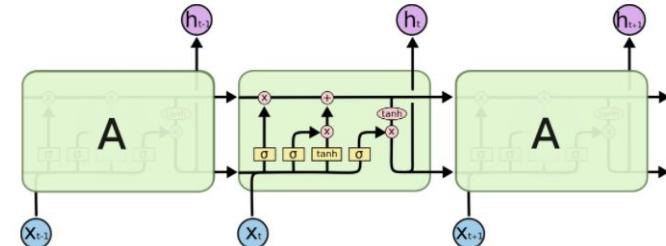


GRU (Gated Recurrent Unit)



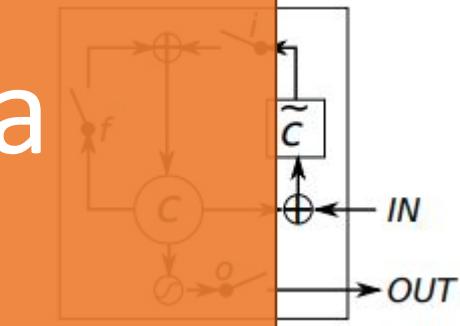
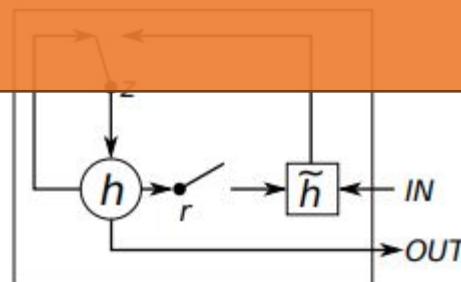
75 other variations

LSTM (Long Short Term Memory)



Same basic idea

75 other variations



Add more neural networks!

- Add “gates” that decide how to balance impact of new input vs hidden state
 - “Remember” gates
 - “Forget” gates
 - “Reset” gates
- Take input and previous state
- Hit ‘em with some weights
- Run it through a sigmoid
- Multiply against h_i or x_i
 - Sigmoid values from 0 – 1 to signify how much to keep/forget

Key Lessons for ML [Domingos, 2012]

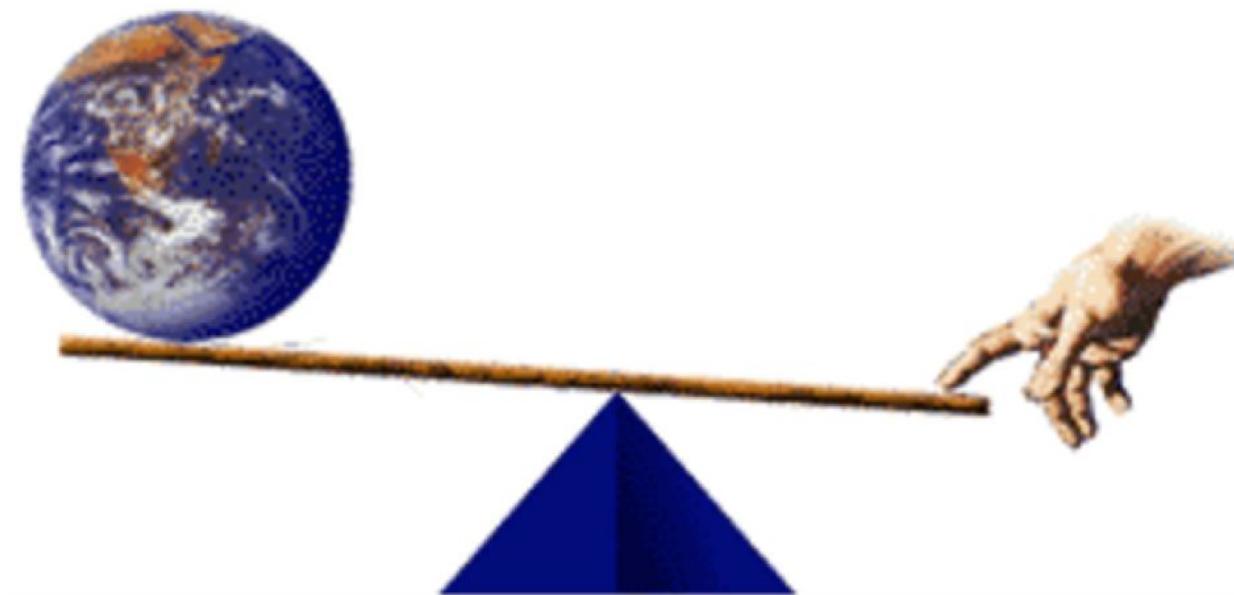
- Learning = Representation + Evaluation + Optimization
- **It's generalization that counts: generalize beyond training examples**
- Data alone is not enough: “no free lunch” theorem--No learner can beat random guessing over all possible functions to be learned
- Intuition fails in high dimensions: “curse 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

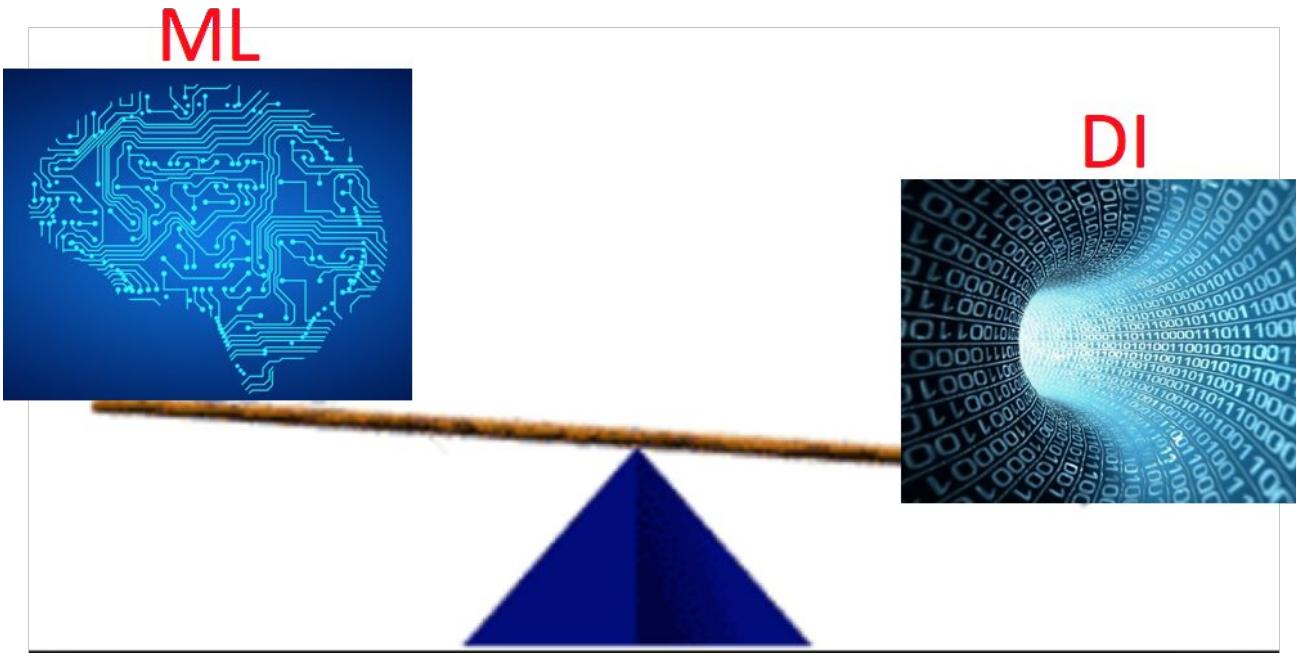
- **ML for effective DI: AUTOMATION, AUTOMATION, AUTOMATION**
 - Automating DI tasks with training data
 - Better understanding of semantics by neural network
- **DI for effective ML: DATA, DATA, DATA**
 - Create large-scale training datasets from different sources
 - Cleaning of data used for training

Give me a Fulscrum, I will Move the Earth

-- Archimedes



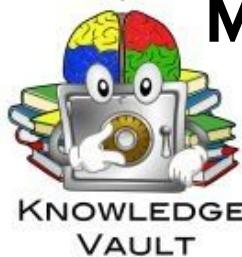
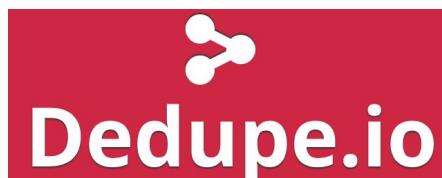
Give me a DI funnel, I will Move ML



Many Systems Where DI & ML Leverage Each Other

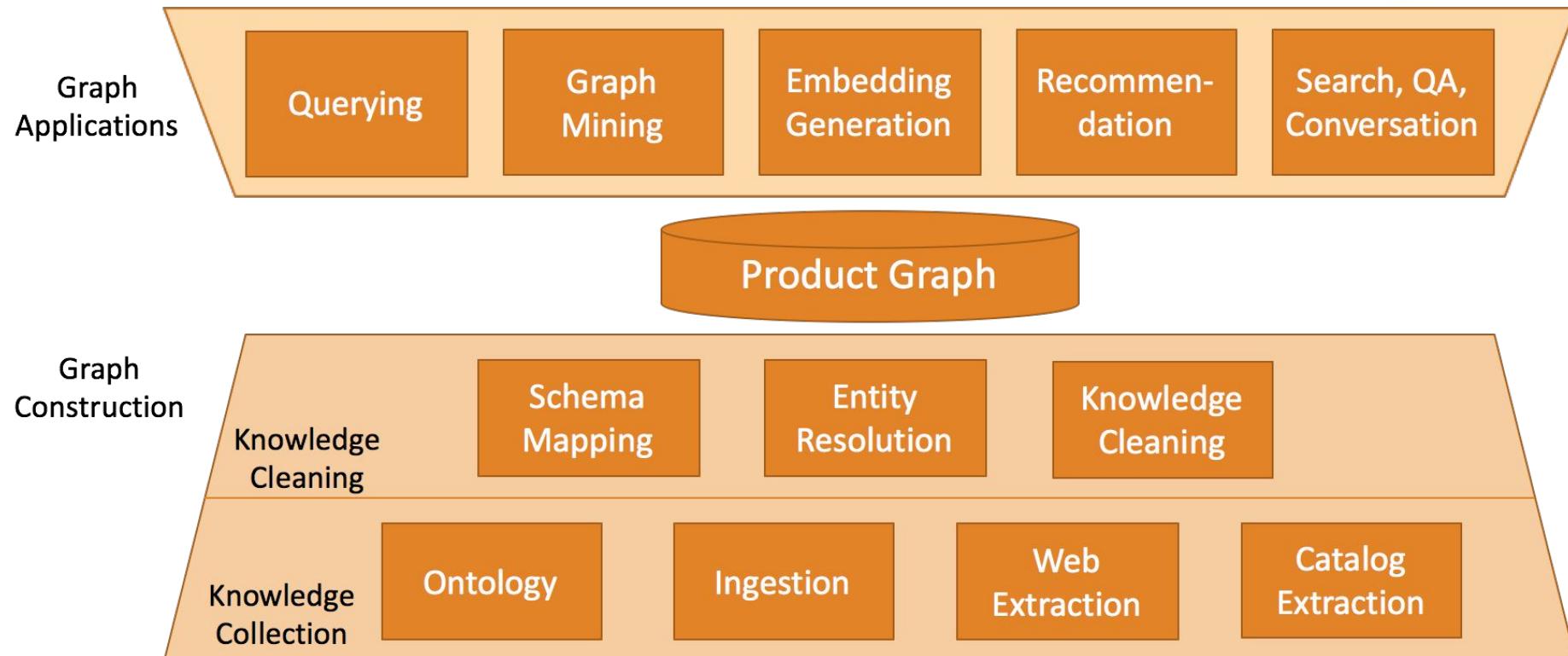


NELL



Increasing number of systems both in industry
and academia.

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



Goal of This Tutorial

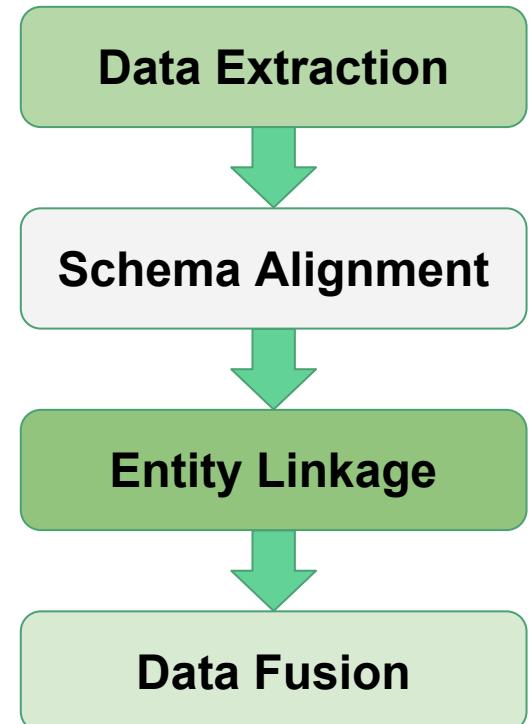
- **NO-GOALS**
 - Present a comprehensive literature review for all topics we are covering
- **GOALS**
 - Present state-of-the-art for DI & ML synergy
 - Show how ML has been transforming DI and vice versa
 - Give some taste on which tool is working best for which tasks
 - Discuss what remains challenging

Outline

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

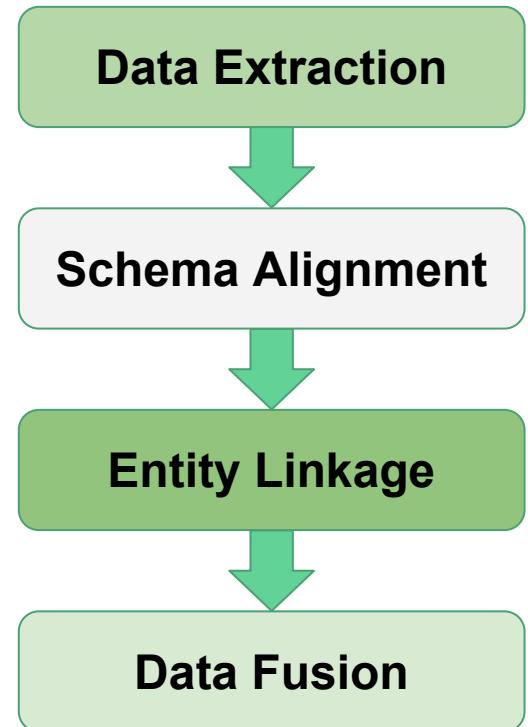
Data Integration Overview

- 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) ³	1789	DOM	Born
2. J. Adams (F)	1797		
3. Jefferson (DR)	1801		
4. Madison (DR)	1809		

Free texts

Synopsis

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.

yelp

DOM

Shana Thai Restaurant

Price Range: \$2

Address: 1000 Market St, San Francisco, CA 94103

Phone: +1 415 981 0000

Website: <http://www.shanathai.com>

Explore the menu

Hours: Monday: 11 am - 1 pm
Tuesday: 11 am - 1 pm
Wednesday: 11 am - 1 pm
Thursday: 11 am - 1 pm
Friday: 11 am - 1 pm
Saturday: 11 am - 1 pm
Sunday: 11 am - 1 pm

Price Range: \$2

Address: 1000 Market St, San Francisco, CA 94103

Phone: +1 415 981 0000

Website: <http://www.shanathai.com>

Explore the menu

Hours: Monday: 11 am - 1 pm
Tuesday: 11 am - 1 pm
Wednesday: 11 am - 1 pm
Thursday: 11 am - 1 pm
Friday: 11 am - 1 pm
Saturday: 11 am - 1 pm
Sunday: 11 am - 1 pm

SCENE FROM "DAN'L DRUCE."

This interesting domestic drama, by Mr. W. S. Gilbert, has continued to engage the sympathies of a mighty sufficient audience at the Lyceum and Theatre Royal, where it has been presented more than a dozen times. Its subject and character were described by us, in the ordinary report of theatrical news, as follows: "The story, as such, need not need to be reminded that the hero of the story, Dan'l Druce, the blacksmith, is a sturdy, manly, drawing on his coat-sleeve, when his coat is off, a valiant fighter from party to party during the civil war of the Commonwealth. His honest, simple, kindly ways, and his love of teetotalism, a helpless female infant is left by some mysterious agency, and may be accepted, as in George Eliot's 'Silas Marner,' as a picture of the true soul-heated misanthrope, far better than riches. In this spirit, at least, he is content to receive the precious human character; and it is this which won him back 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.



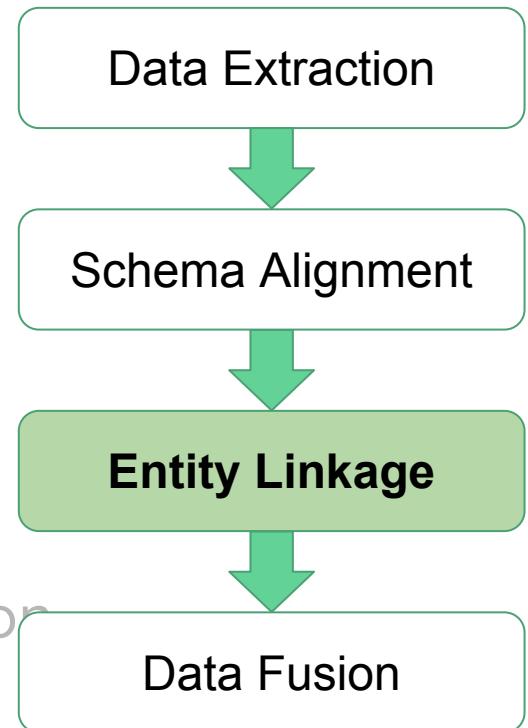
Neural network

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 \mathcal{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

SEE RANK

Anahí 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. Anahí 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 (Q169461)

Mexican singer-songwriter and actress

Mia

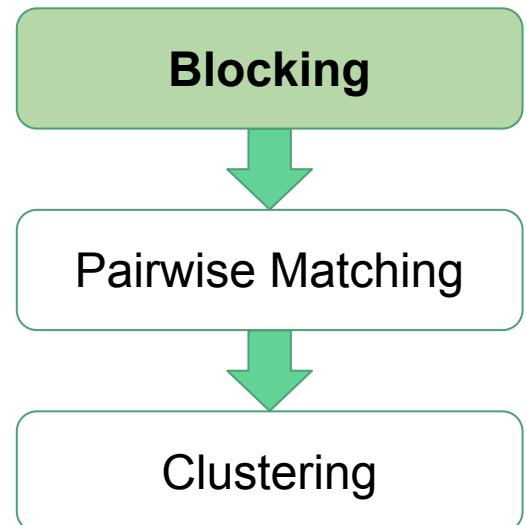
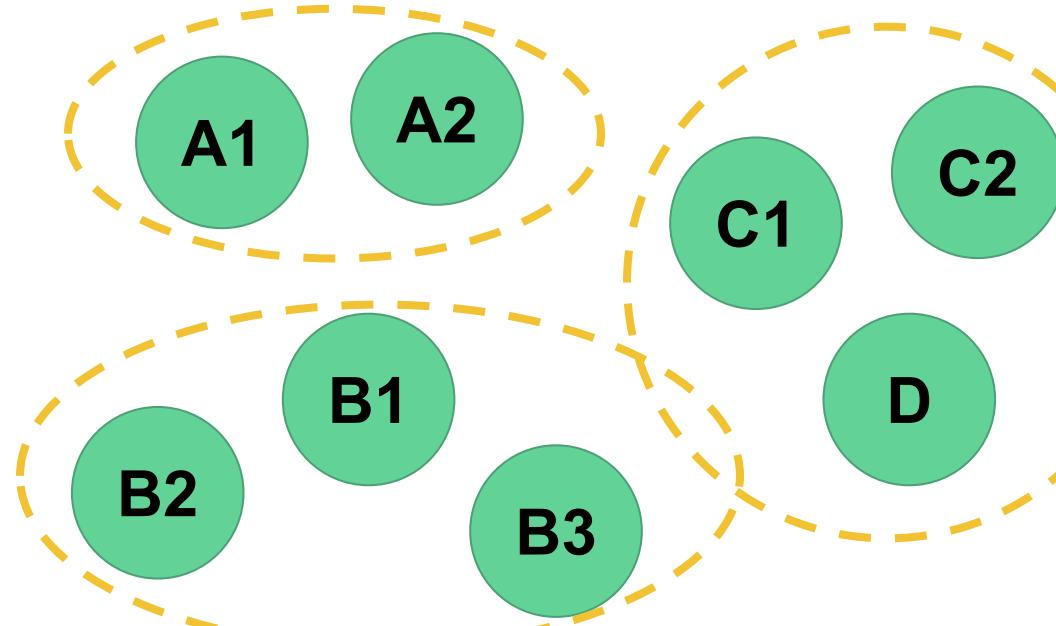
▼ In more languages [Configure](#)

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	edit
	▼ 1 reference	
	imported from	Italian Wikipedia
		+ add reference
		+ add value

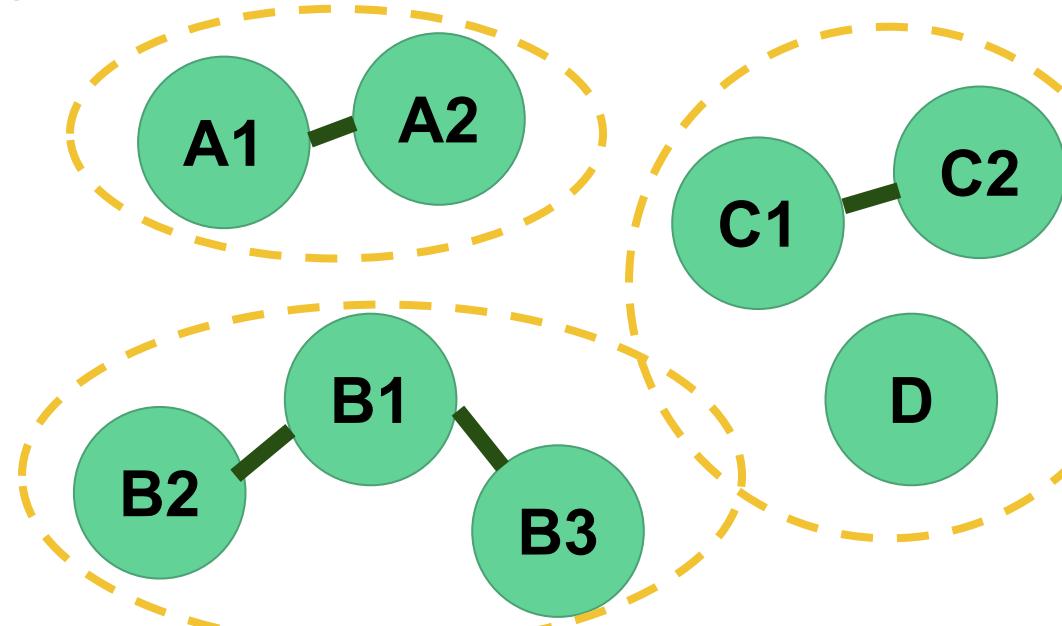
Quick Tour for Entity Linkage

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



Quick Tour for Entity Linkage

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



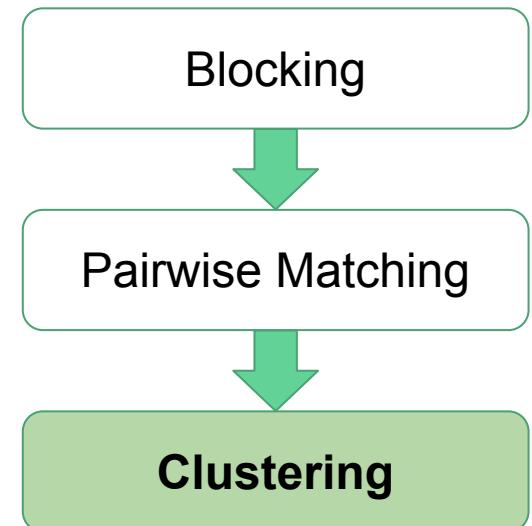
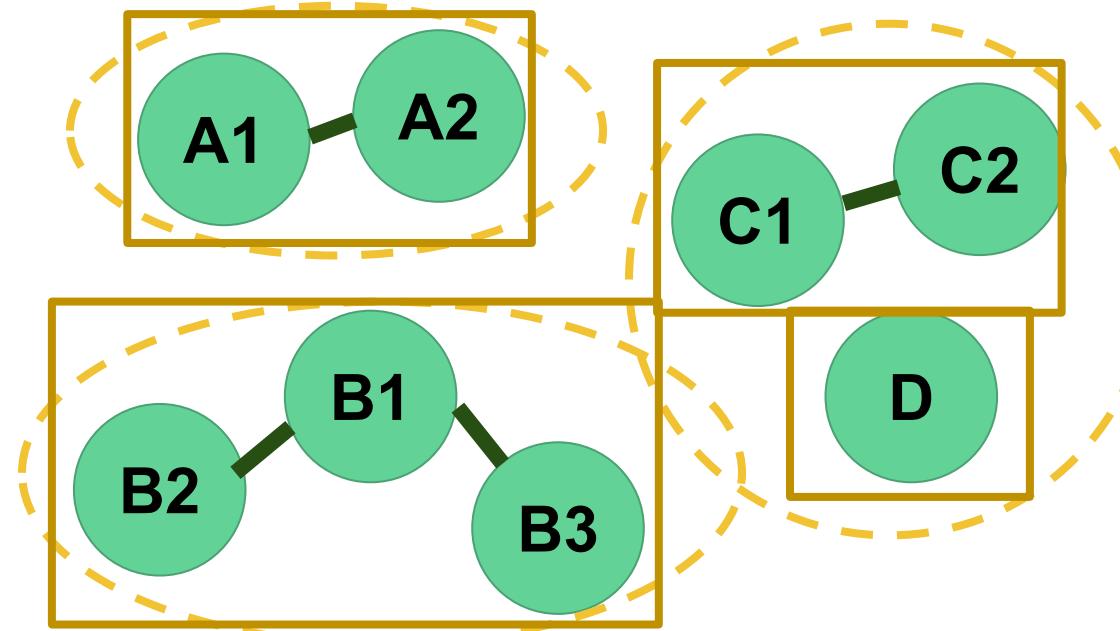
Blocking

Pairwise Matching

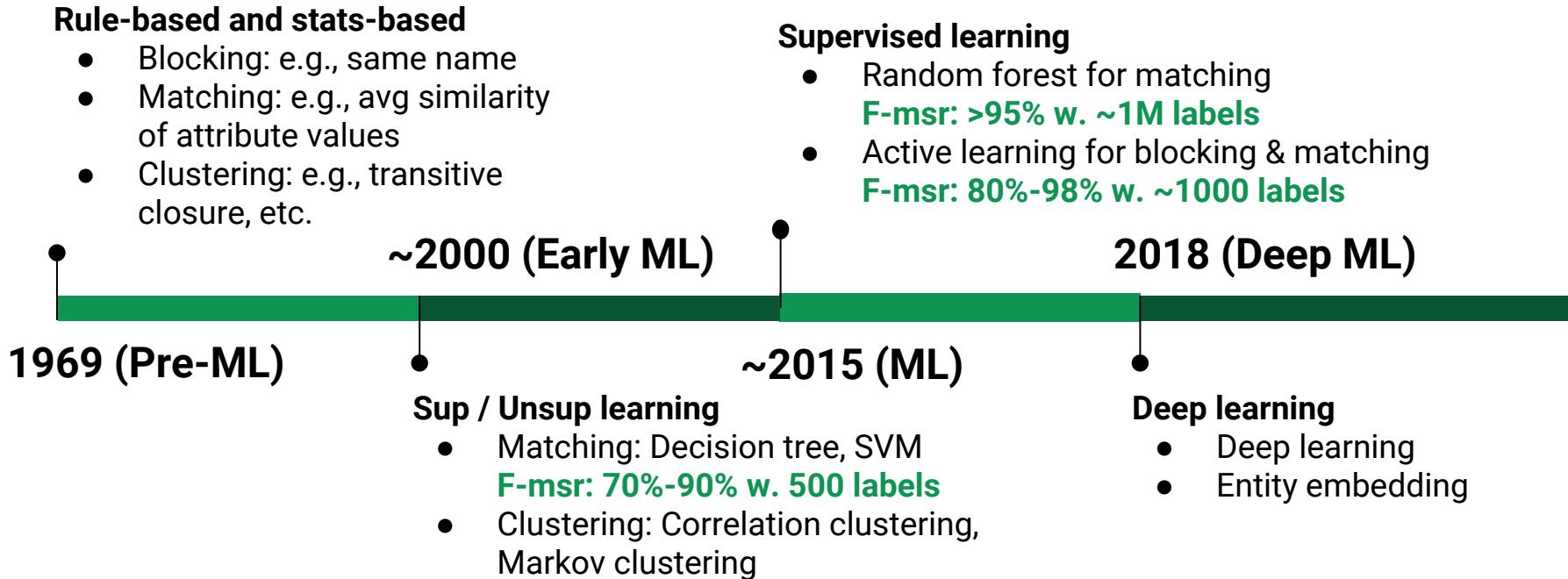
Clustering

Quick Tour for 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.



1969 (Pre-ML)

- [Fellegi and Sunter, 1969]
 - Match: $\text{sim}(r, r') > \theta_h$
 - Unmatch: $\text{sim}(r, r') < \theta_l$
 - Possible match:
$$\theta_l < \text{sim}(r, r') < \theta_h$$

Early ML Models

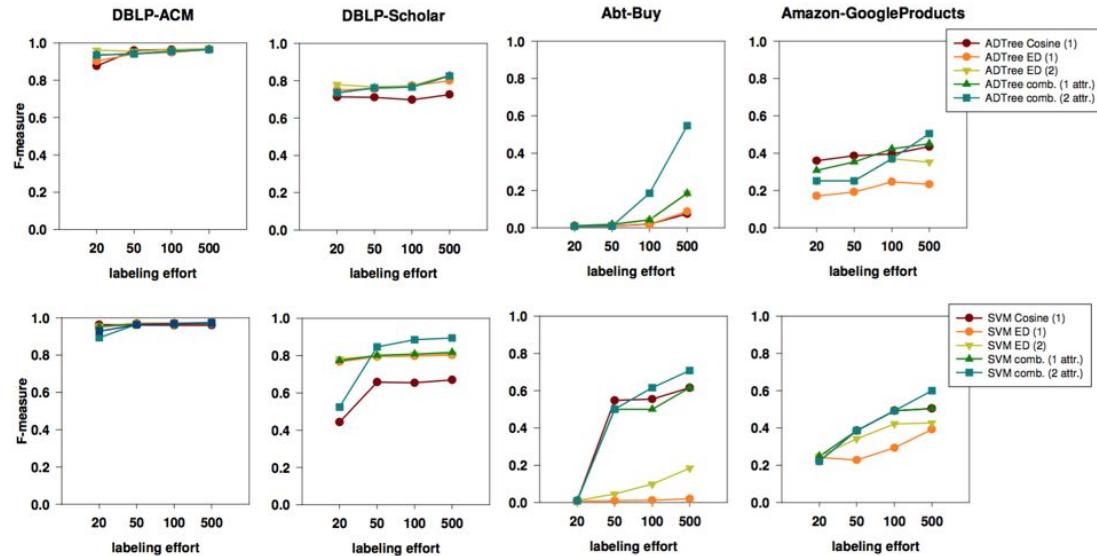
~2000 (Early ML)



Sup / Unsup learning

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

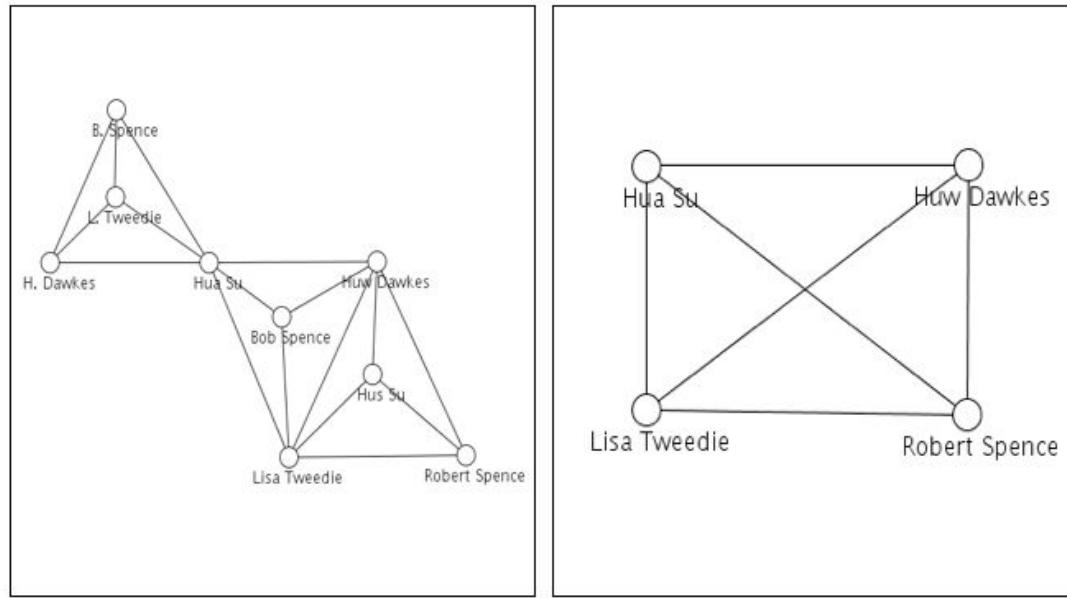
- [Köpcke et al, VLDB'10]



Collective Entity Resolution: Beyond Pairs

- Collective reasoning across entities.
- Constraints across entities:
 - Aggregate constraints
 - Transitivity, Exclusivity
 - Functional dependencies
- Use of probabilistic graphical models, PSL, MLN, to capture such domain knowledge

Out of the scope of this tutorial. For details: See tutorial by Getoor and Machanavajjhala, KDD, 2013.



before

after

[Example by Getoor and Machanavajjhala]

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

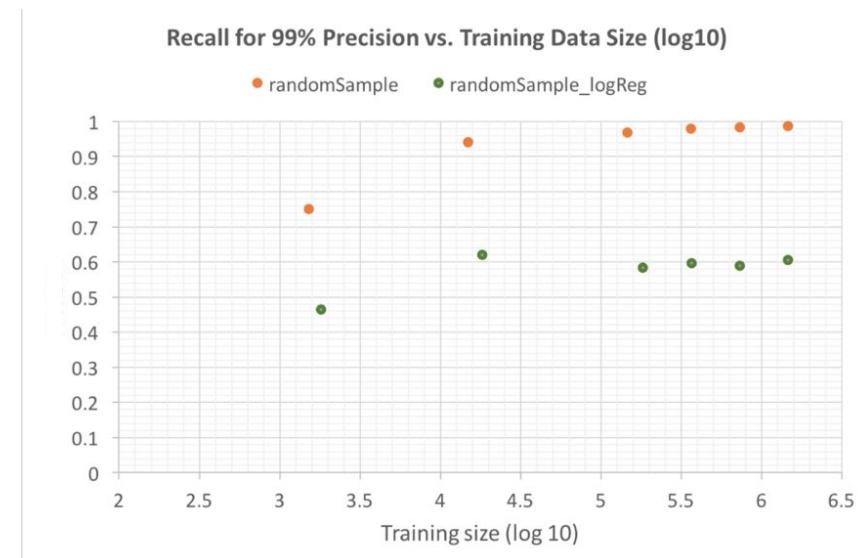
- Random forest for matching
F-msr: >95% w. ~1M labels
- AL for blocking & matching
F-msr: 80%-98% w. ~1000 labels



~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]

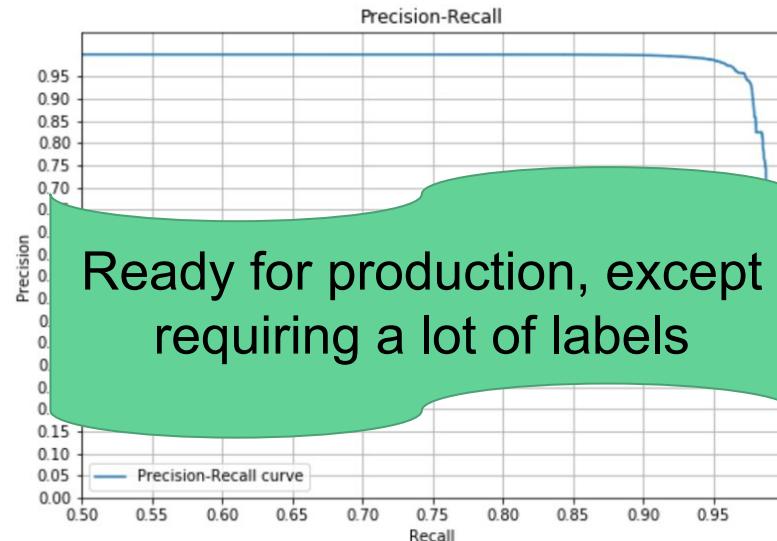
Supervised learning

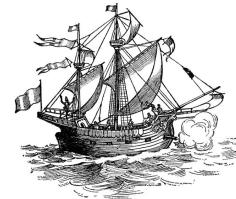
- Random forest for matching
F-msr: >95% w. ~1M labels
- AL for blocking & matching
F-msr: 80%-98% w. ~1000 labels



~2015 (ML)

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





State-of-the-Art ML Models

[Das et al., SIGMOD'17]

Magellan

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)

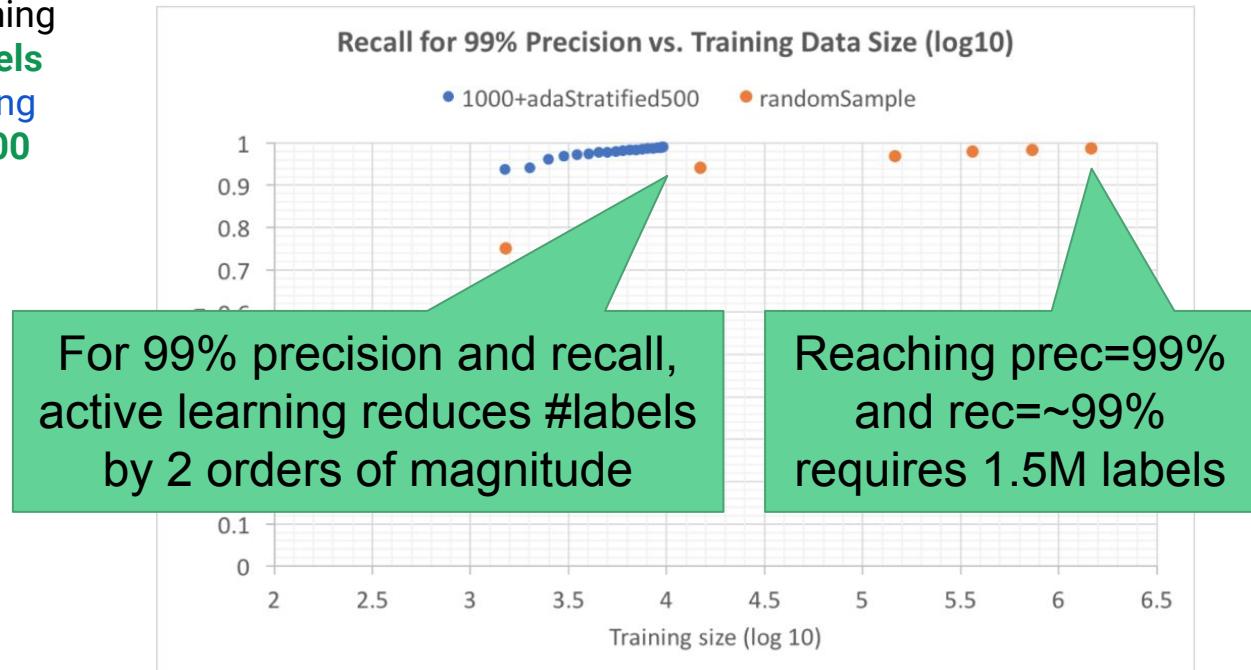
State-of-the-Art ML Models [Dong, KDD'18]

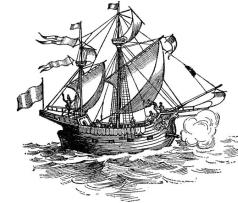
Supervised learning

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



- Apply active learning to minimize #labels





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

Magellan

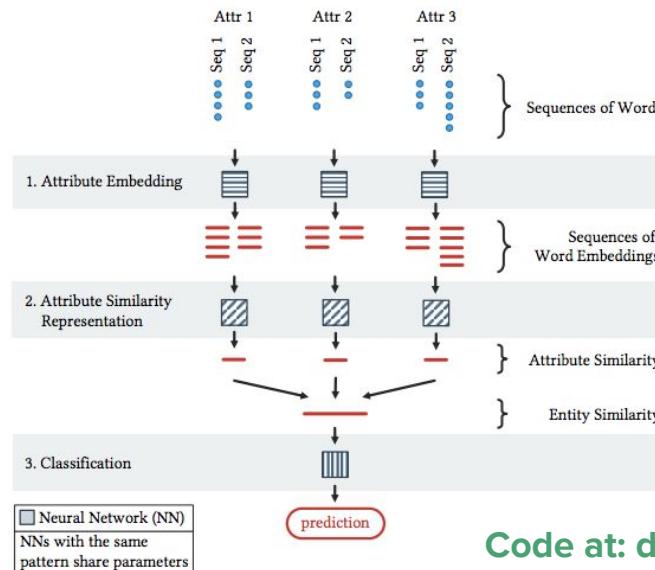
- Embedding on similarities
- Similar performance for structured data;
Significant improvement on texts and dirty data

2018 (Deep ML)



Deep learning

- Deep learning
- Entity embedding



Code at: deepmatcher.ml

Deep Learning Models [Ebraheem et al., VLDB'18]

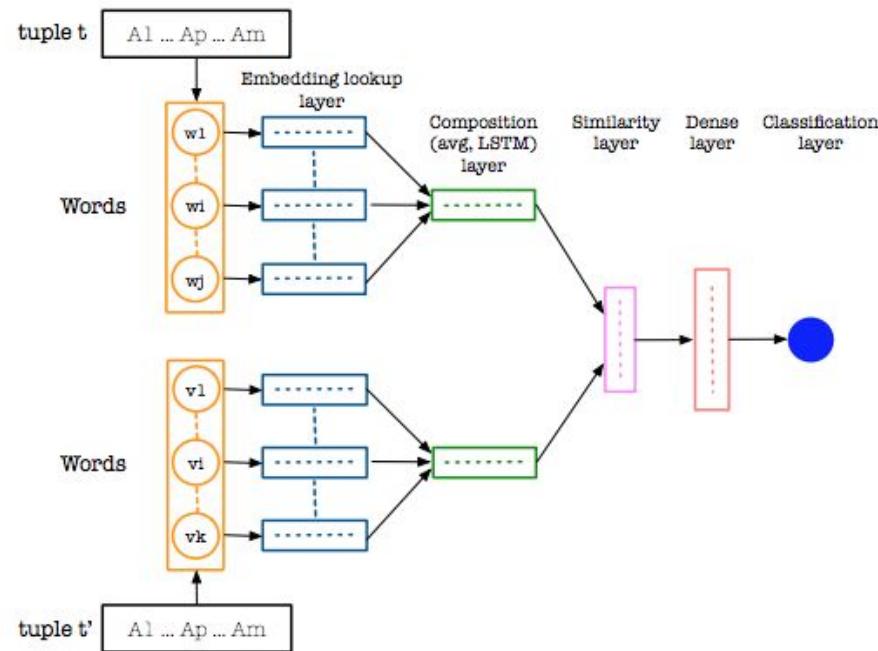
- Embedding on entities
- Outperforming existing solution

2018 (Deep ML)



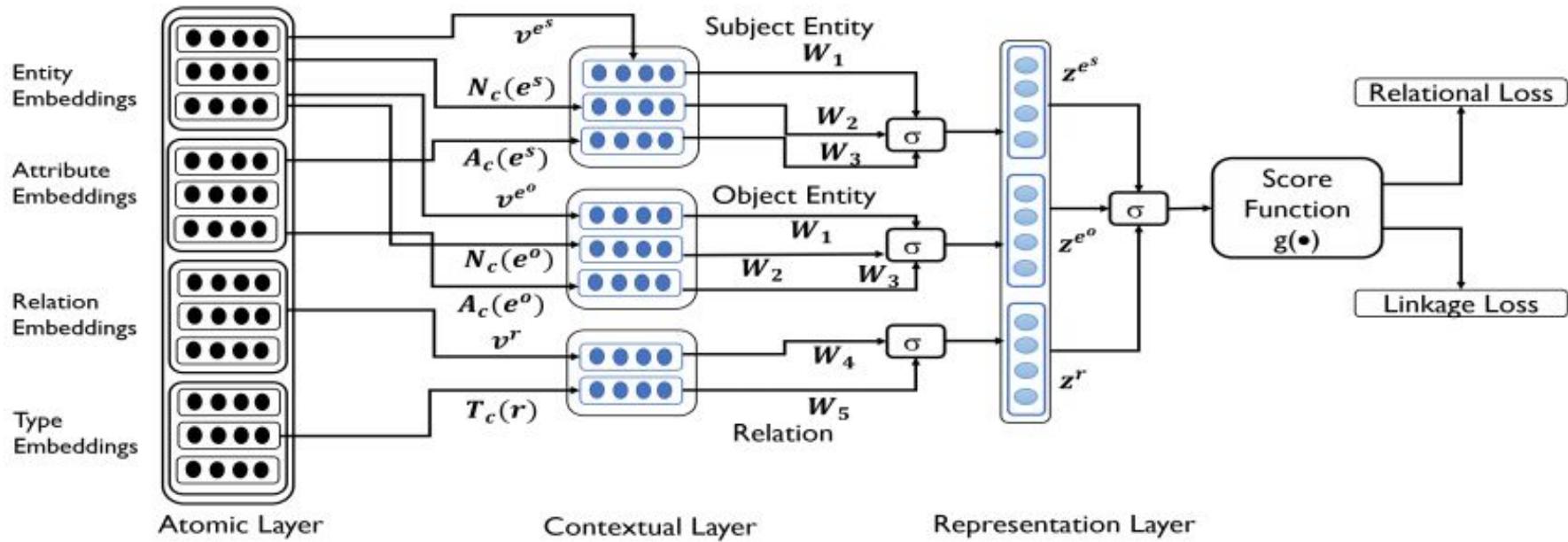
Deep learning

- Deep learning
- Entity embedding



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

- LinkNBed: Embeddings for entities as in knowledge embedding



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

- LinkNBed: 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)

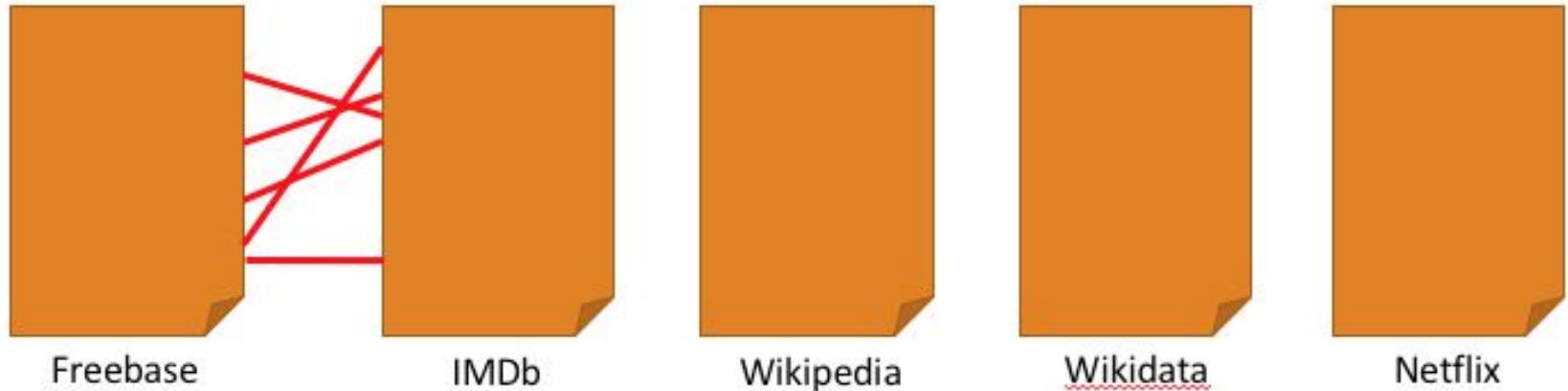


Deep learning

- Deep learning
- Entity embedding

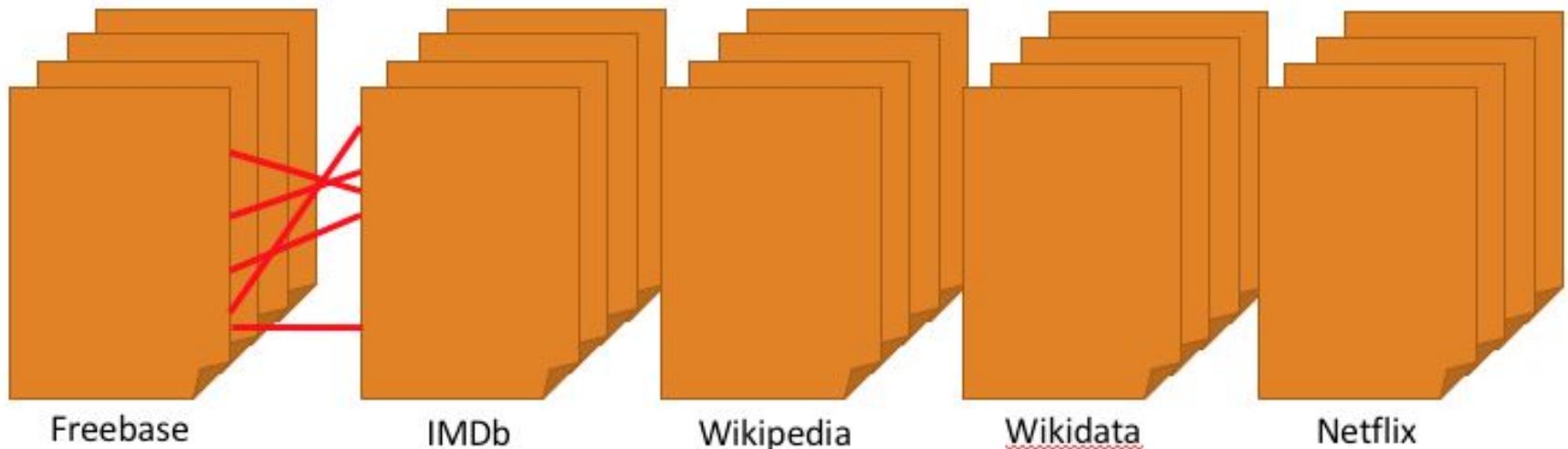
Challenges in Applying ML on EL

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



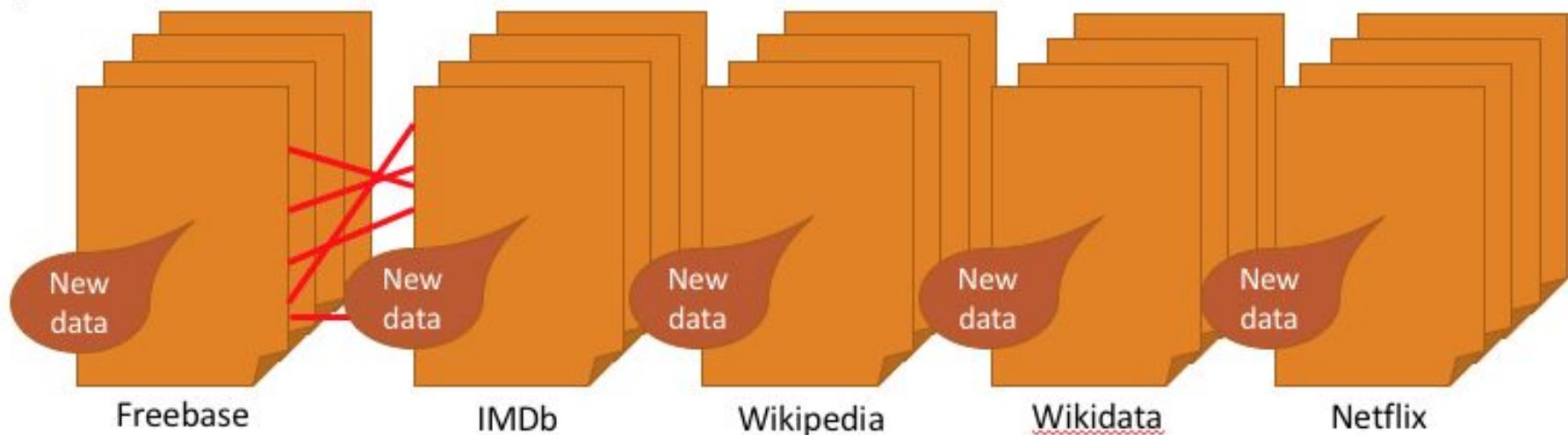
Challenges in Applying ML on EL

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



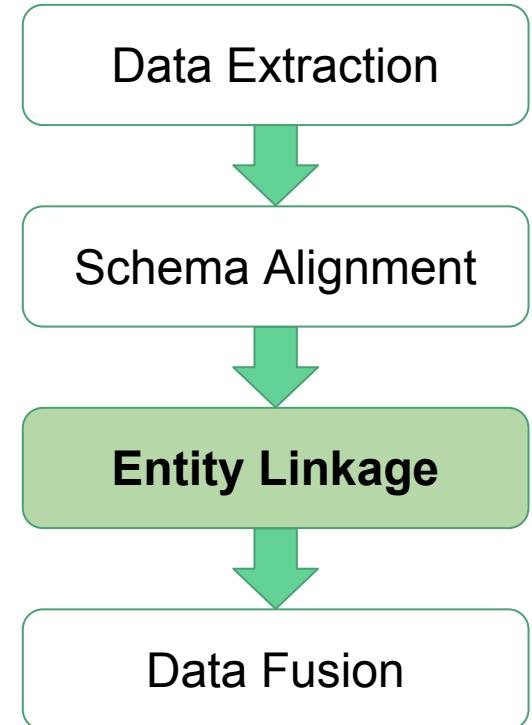
Challenges in Applying ML on EL

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



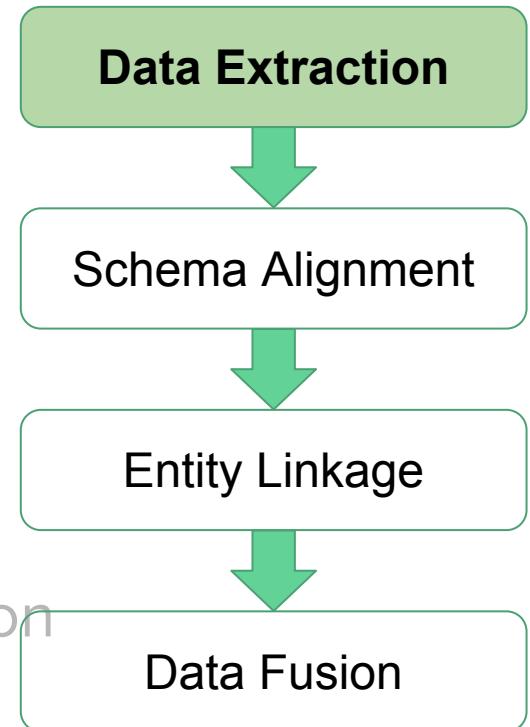
Recipe for Entity Linkage

- 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.

Web tables & Lists

Name and (party) ¹	Term
1. Washington (F) ³	1789–1797
2. J. Adams (F)	1797–1801
3. Jefferson (DR)	1801–1809
4. Madison (DR)	1809–1817

DOM Trees

yelp

Welcome About Me Write a Review Find Friends

Shana Thai Restaurant

Category: Thai (1 star)
311 Moffett Blvd
Bld A
Mountain View, CA 94031
(850) 940-9999
<http://www.shanathai.com>

Explore the menu

Hours:

Price Range: \$

Takes Reservations: No

Delivery: No

Take-out: Yes

Waiter Service: Yes

Outdoor Seating: No

Wi-Fi: No

Good For: Dinner

Free texts

Synopsis

Born on April 15, 1452, in Vinci, Italy, Leonardo da Vinci was a man concerned with the laws of science and nature. He informed his work as a painter, sculptor, architect, engineer, and scientist. His ideas and body of work -- which include the Vitruvian Man, the Last Supper, Leda and the Swan, and the Vitruvian Horse -- influenced countless artists and made him one of the most famous figures of the Italian Renaissance.

Diagram

Regeneration

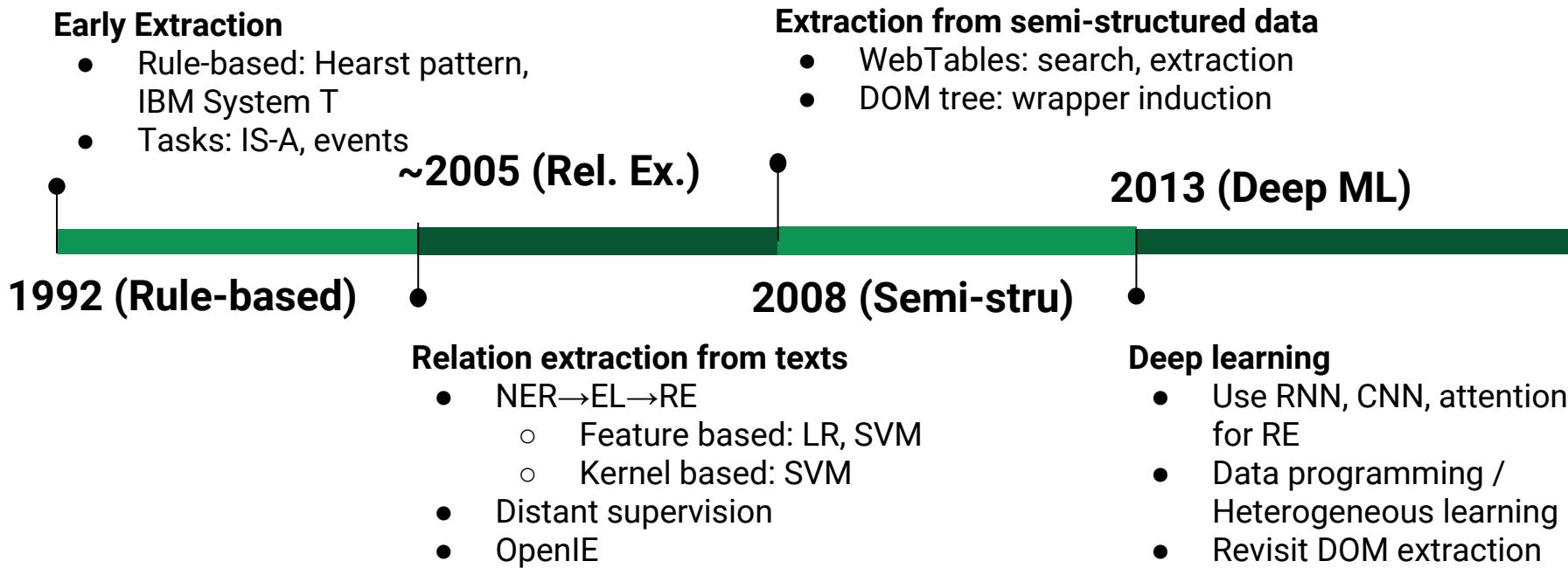
Biological level	Examples	Pre-amputation	Post-amputation	Regenerate
Whole body	Regeneration from a small body fragment			
Structure	Limb, fin, tail, head, tentacle, siphon, arm, stalk			
Internal organ	Heart, liver, lens			
Tissue	Epidermis, gut lining			
Cell	Axon, muscle fiber			

TRENDS in Ecology & Evolution

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")

35 Years of Data Extraction



Extraction from Texts: Quick Tour

Bill Gates founded Microsoft in 1975.

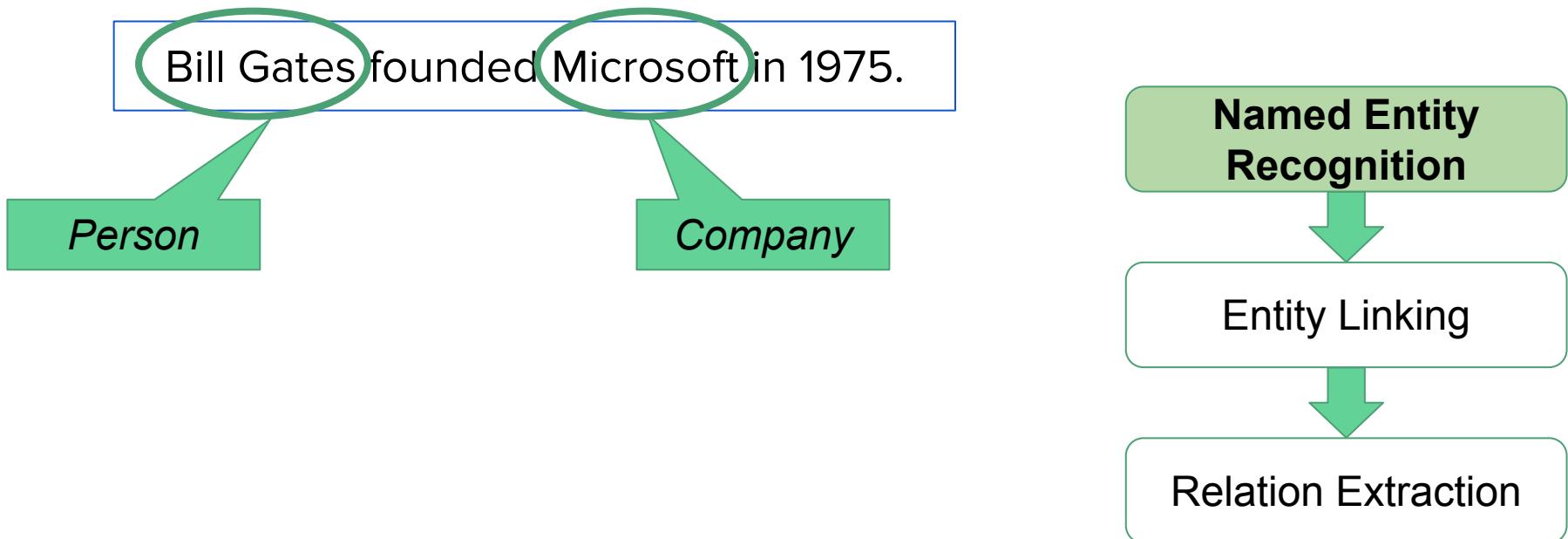
Named Entity
Recognition

Entity Linking

Relation Extraction



Extraction from Texts: Quick Tour



Extraction from Texts: Quick Tour

Bill Gates founded Microsoft in 1975.



Entity **linkage**: linking two structured records
Entity **linking**: linking a phrase in texts to an entity in a reference list (e.g., knowledge graph)

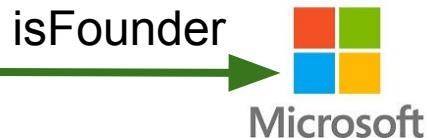
Named Entity Recognition

Entity Linking

Relation Extraction

Extraction from Texts: Quick Tour

Bill Gates founded Microsoft in 1975.



We focus on Relation Extraction in the rest of the tutorial.

Named Entity Recognition

Entity Linking

Relation Extraction

Extraction from Texts: Feature Based [Zhou et al., ACL'05]

~2005 (Rel. Ex.)



Relation extraction from texts

- NER→EL→RE
 - Feature based: LR, SVM
 - Kernel based: SVM
- Distant supervision
- OpenIE

- **Models**
 - Logistic regression
 - SVM (Support Vector Machine)
- **Features**
 - Lexical: entity, part-of-speech, neighbor
 - Syntactic: **chunking**, parse tree
 - Semantic: concept hierarchy, entity class
- **Results**
 - Prec=≈60%, Rec=≈50%

Extraction from Texts: Feature Based [Zhou et al., ACL'05]

~2005 (Rel. Ex.)



Relation extraction from texts

- NER→EL→RE
 - Feature based: LR, SVM
 - Kernel based: SVM
- Distant supervision
- OpenIE

Features	P	R	F
Words	69.2	23.7	35.3
+Entity Type	67.1	32.1	43.4
+Mention Level	67.1	33.0	44.2
+Overlap	57.4	40.9	47.8
+Chunking	61.5	46.5	53.0
+Dependency Tree	62.1	47.2	53.6
+Parse Tree	62.3	47.6	54.0
+Semantic Resources	63.1	49.5	55.5

Major Lift

Table 2: Contribution of different features over 43 relation subtypes in the test data

Extraction from Texts: Kernel Based [Mengqiu Wang, IJCNLP'08]

~2005 (Rel. Ex.)



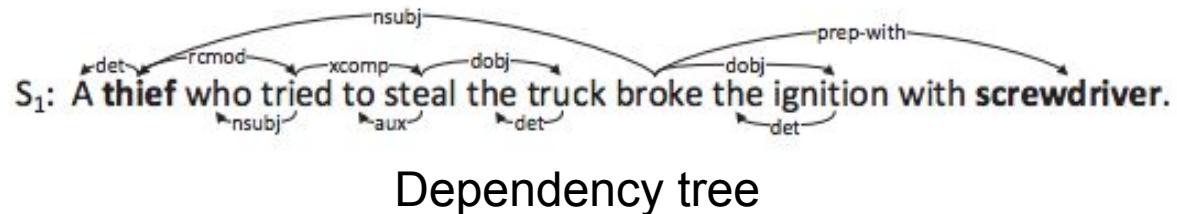
Relation extraction from texts

- NER→EL→RE
 - Feature based: LR, SVM
 - Kernel based: SVM
- Distant supervision
- OpenIE

- **Models**
 - SVM (Support Vector Machine)
- **Kernels**
 - Subsequence
 - Dependency tree
 - **Shortest dependency path**
 - Convolution dependency

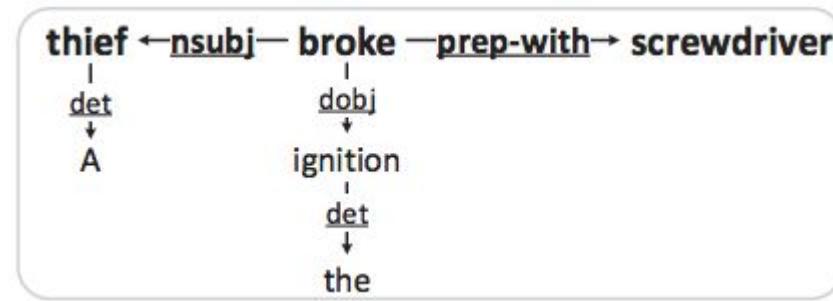
Extraction from Texts: Kernel Based [Mengqiu Wang, IJCNLP'08]

~2005 (Rel. Ex.)



Relation extraction from texts

- NER→EL→RE
 - Feature based: LR, SVM
 - Kernel based: SVM
- Distant supervision
- OpenIE



Extraction from Texts: Kernel Based [Mengqiu Wang, IJCNLP'08]

~2005 (Rel. Ex.)



Relation extraction from texts

- NER→EL→RE
 - Feature based: LR, SVM
 - Kernel based: SVM
- Distant supervision
- OpenIE

- Models
 - SVM (Support Vector Machine)
- Kernels
 - Subsequence
 - Dependency tree
 - Shortest dependency path
 - Convolution dependency
- Results
 - Prec=≈70%, Rec=≈40%

Extraction from Texts: Kernel Based [Mengqiu Wang, IJCNLP'08]

~2005 (Rel. Ex.)



Relation extraction from texts

- NER→EL→RE
 - Feature based: LR, SVM
 - Kernel based: SVM
- Distant supervision
- OpenIE

kernel method	5-fold CV on ACE 2003		
	Precision	Recall	F1
subsequence	0.703	0.389	0.546
dependency tree	0.681	0.290	0.485
shortest path	0.747	0.376	0.562

Table 1: Results of different kernels on ACE 2003 training set using 5-fold cross-validation.

Extraction from Texts: Deep Learning

2013 (Deep ML)

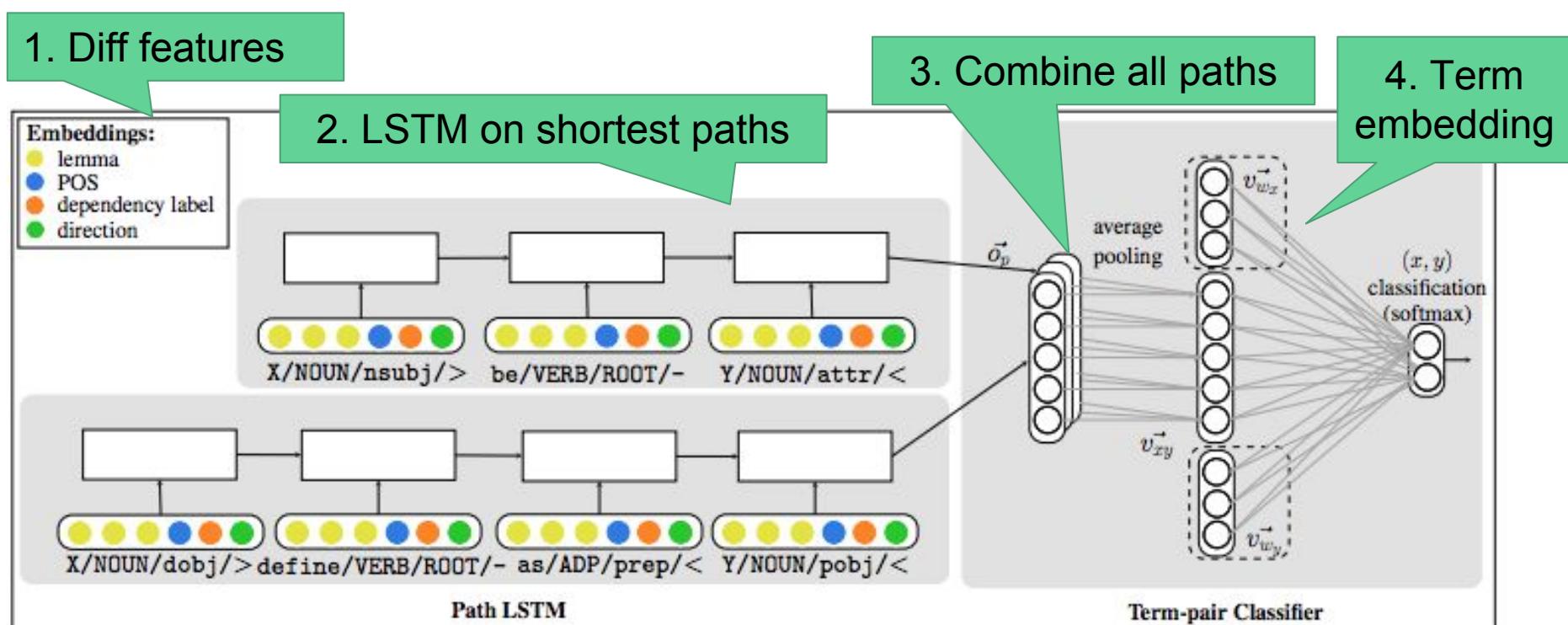


Deep learning

- Use RNN, CNN, attention for RE
- Data programming / Heterogeneous learning
- Revisit DOM extraction

- Same intuitions, different models
 - (2012-13) Recursive NN: dependency tree [Socher et al., EMNLP'12] [Hashimoto et al., EMNLP'13]
 - (2014-15) CNN: shortest dependency path [Zeng et al., COLING'14][Liu et al., ACL '15]
 - (2015+) LSTM: shortest dependency path, lexical/syntactic/semantic features [Xu et al., EMNLP'15][Shwartz et al., ACL'16] [Nguyen, NAACL'16]

Example System: HyperNET [Shwartz et al., ACL'16]



Quality in identifying hypernyms: Prec = 0.9, Rec = 0.9

Label Generation for Extraction Training

Where are training labels from?

~2005 (Rel. Ex.)



- **Semi-supervised learning**
 - Iterative extraction [Carlson et al., AAAI'10]
Use new extractions to retrain models
E.g., NELL

Relation extraction from texts

- **NER→EL→RE**
 - Feature based: LR, SVM
 - Kernel based: SVM
- Distant supervision
- OpenIE

Iterations	Estimated Precision (%)	# Promotions
1–22	90	88,502
23–44	71	77,835
45–66	57	76,116

Label Generation for Extraction Training

Where are training labels from?

~2005 (Rel. Ex.)



Relation extraction from texts

- NER→EL→RE
 - Feature based: LR, SVM
 - Kernel based: SVM
- Distant supervision
- OpenIE

- Semi-supervised learning
 - Iterative extraction [Carlson et al., AAAI'10]
Use new extractions to retrain models
E.g., NELL
- Weak learning
 - Distant supervision [Mintz et al., ACL'09]
Rule-based annotation with seed data
E.g., DeepDive, Knowledge Vault

Will cover in “DI for ML”

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.

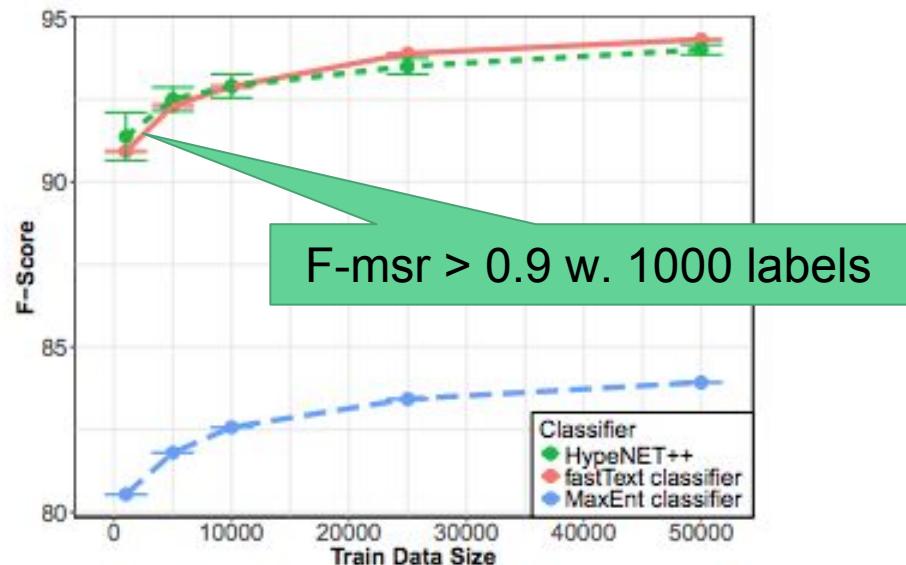
[Adapted example from Luke Zettlemoyer]

Label Generation for Extraction Training

~2005 (Rel. Ex.)

Where are training labels from?

- Distant supervision: HyperNet++
[Christodoulopoulos & Mittal, 18]



Label Generation for Extraction Training

Where are training labels from?

2013 (Deep ML)



Deep learning

- Use RNN, CNN, attention for RE
- Data programming / Heterogeneous learning
- Revisit DOM extraction

Will cover in “DI for ML”

- **Semi-supervised learning**
 - Iterative extraction [Carlson et al., AAAI’10]
Use new extractions to retrain models
E.g., NELL
- **Weak learning**
 - Distant supervision [Mintz et al., ACL’09]
Rule-based annotation with seed data
E.g., DeepDive, Knowledge Vault
 - Data programming [Ratner et al., NIPS’16]
Manually write labelling functions
E.g., Snorkle, Fouduer

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

Input: Labeling Functions,
Unlabeled data

DOMAIN
EXPERT

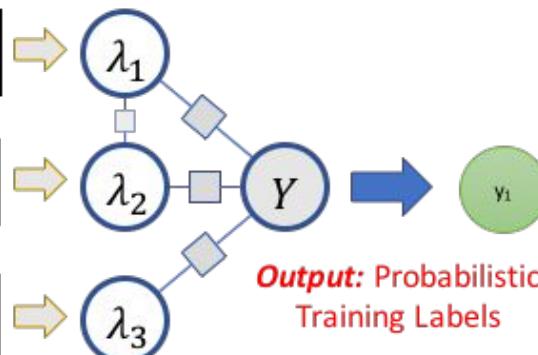


```
def lf1(x):
    cid = (x.chemical_id,
           x.disease_id)
    return 1 if cid in KB else 0
```

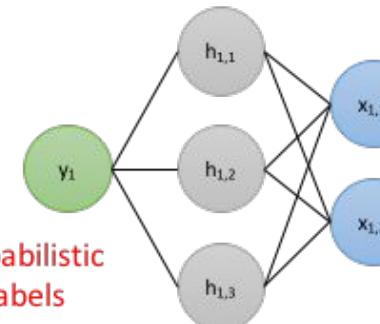
```
def lf2(x):
    m = re.search(r'.*cause.*',
                  x.between)
    return 1 if m else 0
```

```
def lf3(x):
    m = re.search(r'.*not
                  cause.*',
                  x.between)
    return 1 if m else 0
```

**Generative
Model**



**Noise-Aware
Discriminative Model**



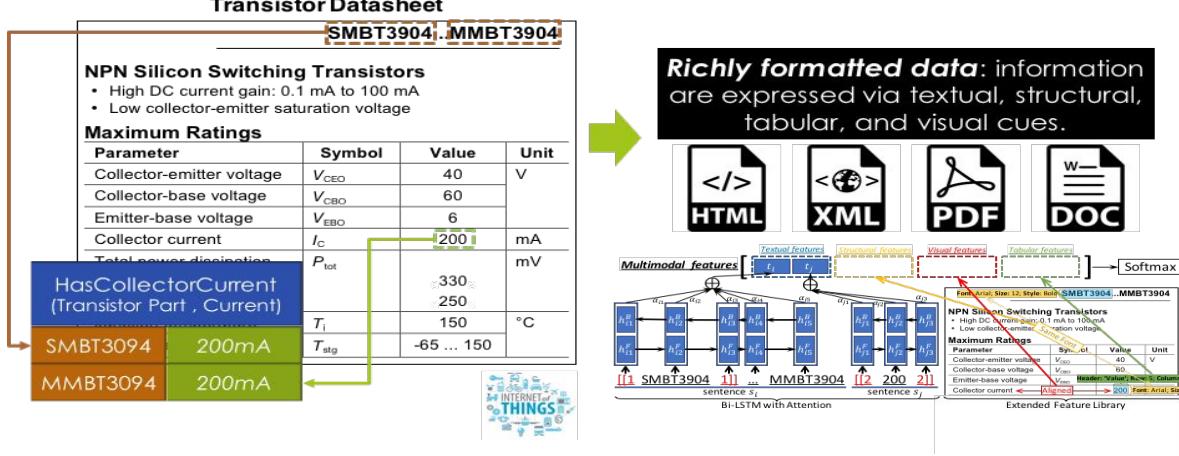
*Ex. Application:
Knowledge Base
Creation (KBC)*



- 1 Users write *labeling functions* to generate noisy labels
- 2 We model the labeling functions' behavior to de-noise them
- 3 We use the resulting prob. labels to train a model



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



Fonduer combines a new **BiLSTM with multimodal features and data programming**.

System	ELEC.	GEN.	
	Digi-Key	GWAS Central	GWAS Catalog
Knowledge Base			
# Entries in KB	376	3,008	4,023
# Entries in Fonduer	447	6,420	6,420
Coverage	0.99	0.82	0.80
Accuracy	0.87	0.87	0.89
# New Correct Entries	17	3,154	2,486
Increase in Correct Entries	1.05×	1.87×	1.42×

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

OpenIE from Texts

~2005 (Rel. Ex.)



Where are predicates from?

- **ClosedIE**
 - Only extracting facts corresponding to ontology
 - Normalize predicates by ontology
 - E.g., (Bill Gates, /person/isFounder, Microsoft)
- **Relation extraction from texts**
 - NER→EL→RE
 - Feature based: LR, SVM
 - Kernel based: SVM
 - Distant supervision
 - **OpenIE**
- **Bill Gates founded Microsoft in 1975.**
- **OpenIE** [Banko et al., IJCAI'07]
 - Extract all relations expressed in texts
 - Predicates are unnormalized strings
 - E.g., (“Bill Gates”, “founded”, “Microsoft”)

OpenIE from Texts [Etzioni et al., IJCAI'11]

ClosedIE

Named Entity
Recognition



Entity Linking



Relation Extraction

OpenIE

Predicate
Identification



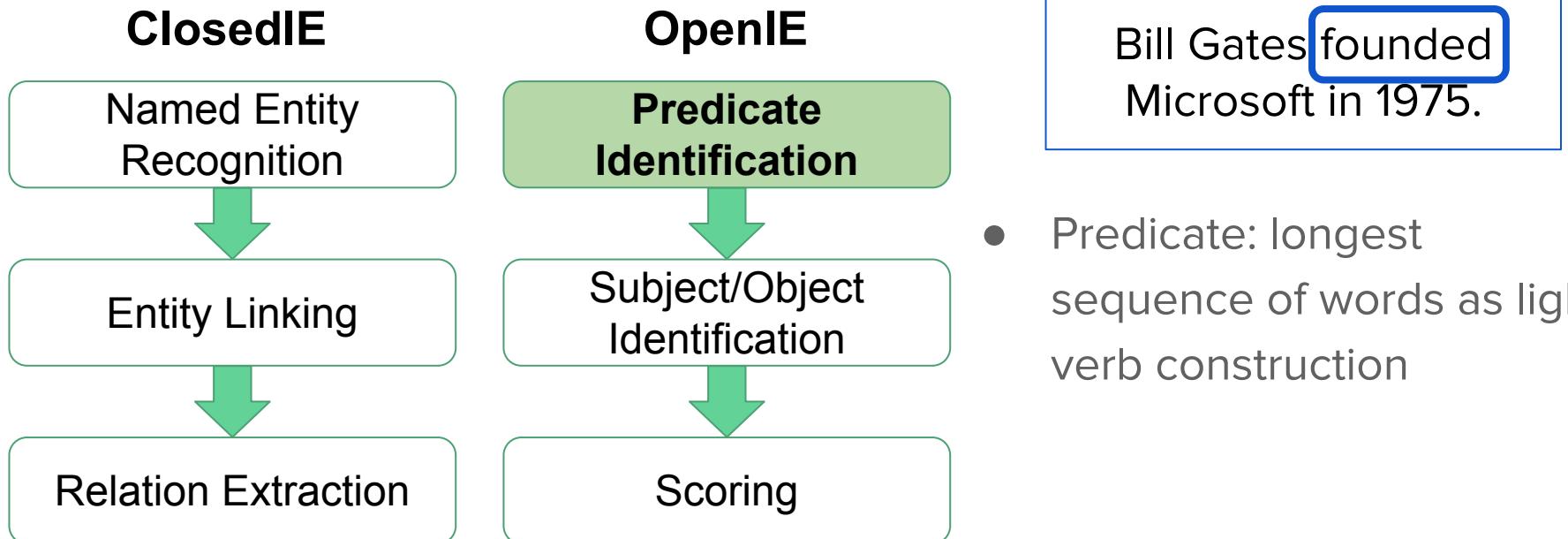
Subject/Object
Identification



Scoring

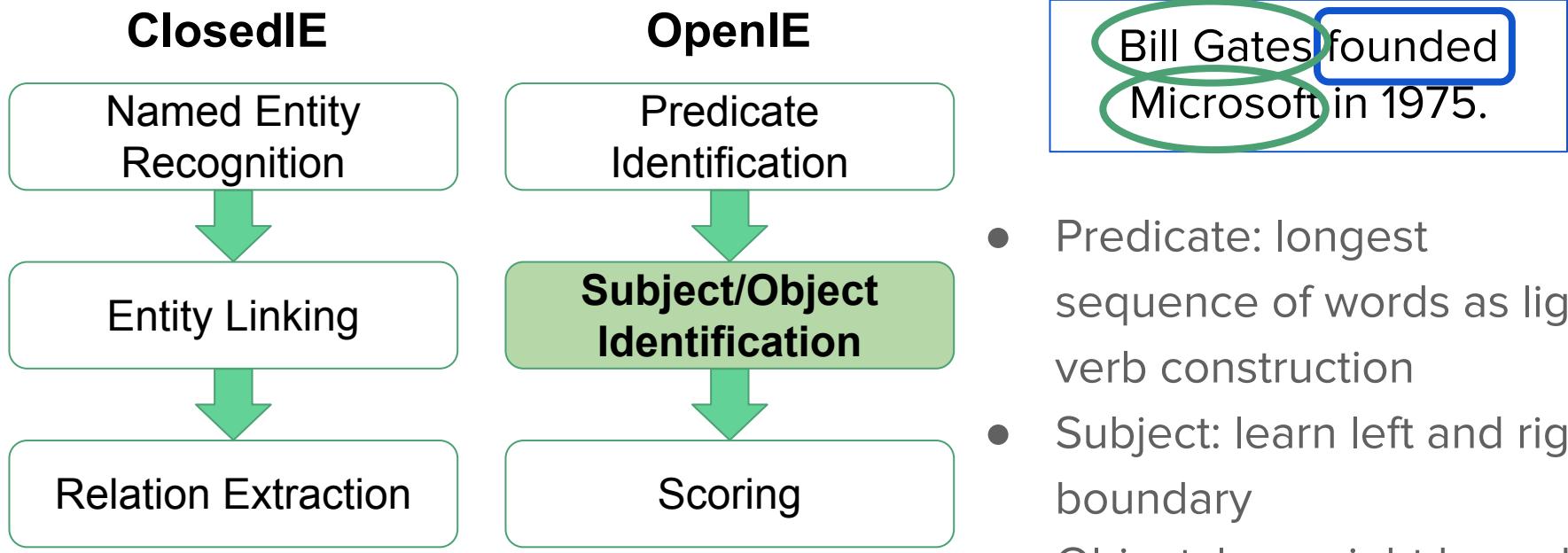
Bill Gates founded
Microsoft in 1975.

OpenIE from Texts [Etzioni et al., IJCAI'11]



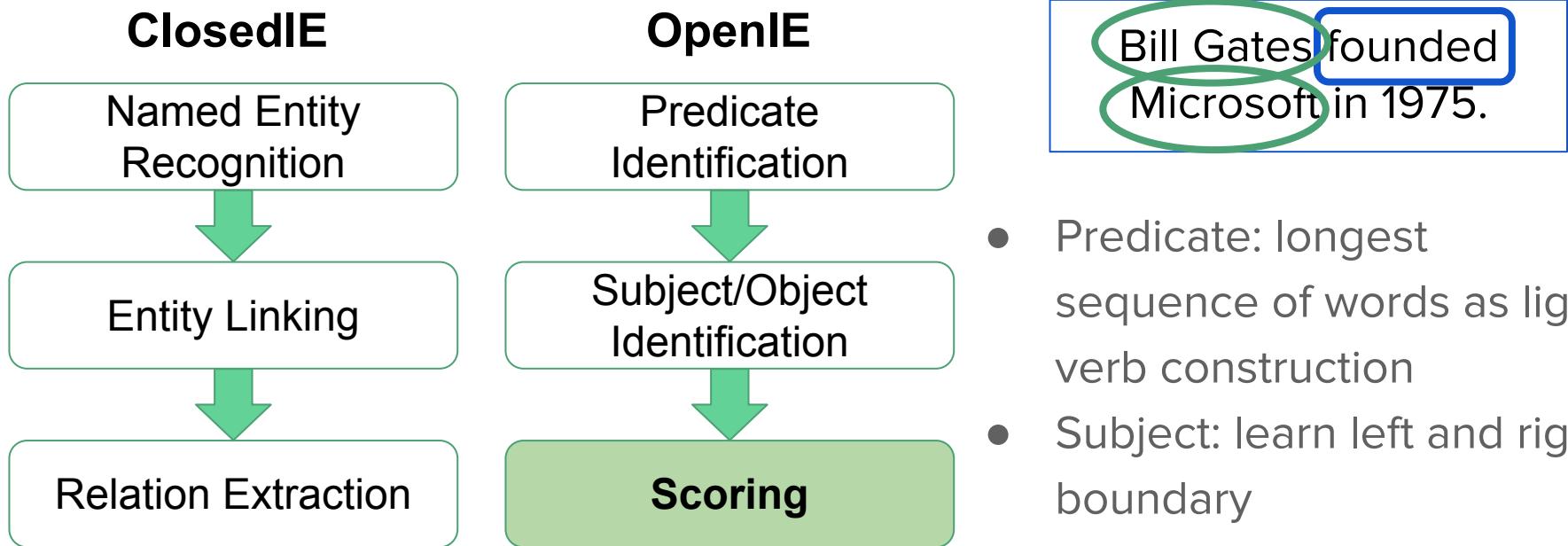
- Predicate: longest sequence of words as light verb construction

OpenIE from Texts [Etzioni et al., IJCAI'11]



- Predicate: longest sequence of words as light verb construction
- Subject: learn left and right boundary
- Object: learn right boundary

OpenIE from Texts [Etzioni et al., IJCAI'11]



- Predicate: longest sequence of words as light verb construction
- Subject: learn left and right boundary
- Object: learn right boundary
- LR for triple confidence

OpenIE from Texts [Mausam et al., EMNLP'12]

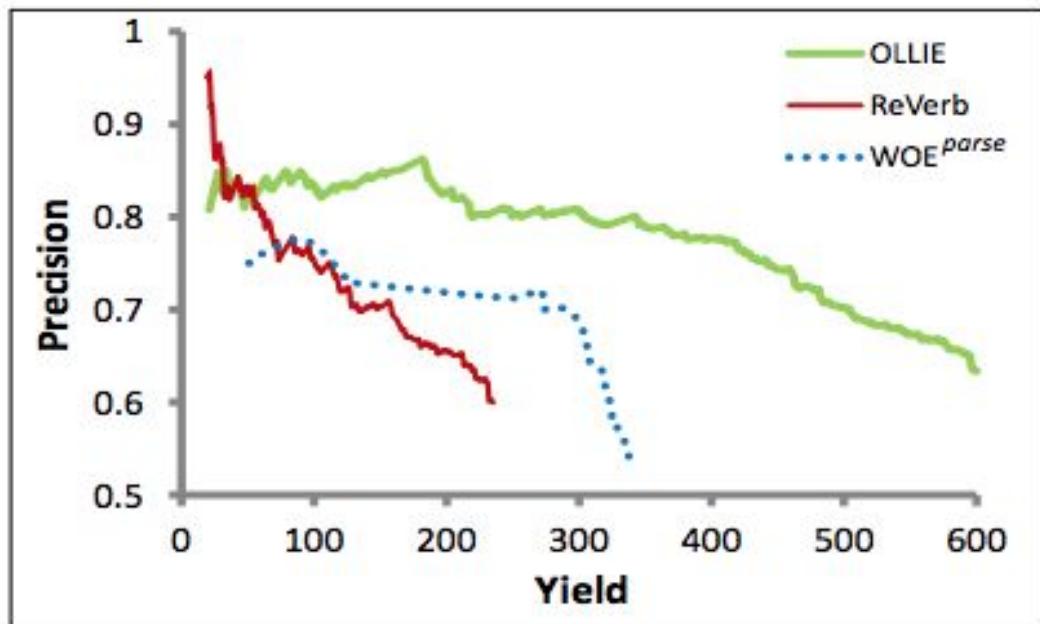
Where are predicates from?

~2005 (Rel. Ex.)



Relation extraction from texts

- NER→EL→RE
 - Feature based: LR, SVM
 - Kernel based: SVM
- Distant supervision
- OpenIE



Extraction from Semi-Structured Data

Extraction from semi-structured data

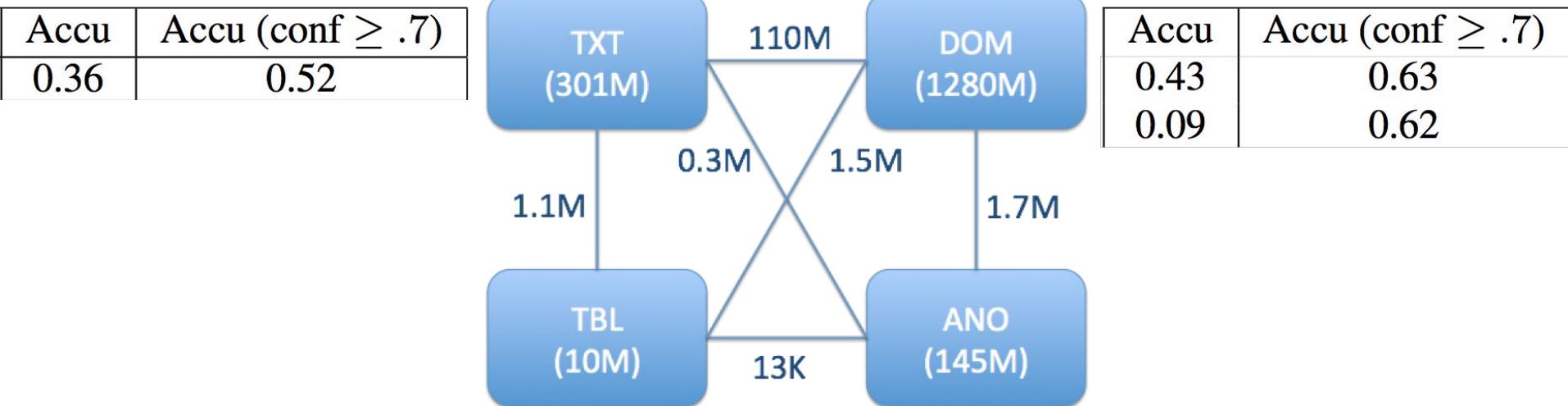
- WebTables: search, extraction
- DOM tree: wrapper induction



2008 (Semi-stru)

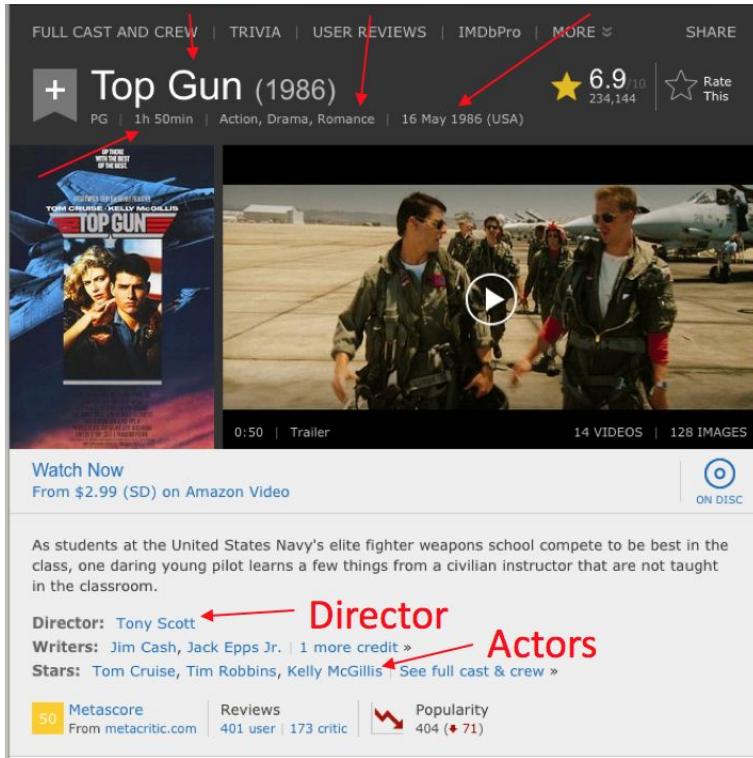
Why Semi-Structured Data?

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



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

Runtime



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

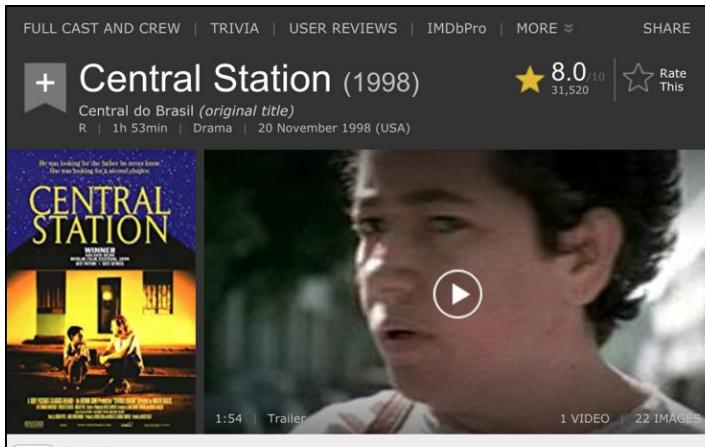
Actor (46 credits) Hide ▲

Top Gun: Maverick (<i>pre-production</i>) Maverick	2019
M:I 6 - Mission Impossible (<i>filming</i>) Ethan Hunt	2018
American Made (<i>completed</i>) Barry Seal	2017
Luna Park (<i>announced</i>)	
The Mummy Nick Morton	2017
Jack Reacher: Never Go Back Jack Reacher	2016
Mission: Impossible - Rogue Nation Ethan Hunt	2015
Edge of Tomorrow Cage	2014
Oblivion Jack	2013/I
Jack Reacher Reacher	2012
Rock of Ages Stacee Jaxx	2012
Mission: Impossible - Ghost Protocol Ethan Hunt	2011
Knight and Day Roy Miller	2010
Valkyrie Colonel Claus von Stauffenberg	2008
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_flm_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/ch0085702/?ref=nm_flm_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_flm_act_5">The Mummy</a>
    </b>
    <br>
      <a href="/character/ch0573416/?ref=nm_flm_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-tti1631867">..</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-tt1336608">..</div>
```

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

- Challenge: slight variations from page to page



FULL CAST AND CREW | TRIVIA | USER REVIEWS | IMDbPro | MORE ▾ SHARE

+ **Central Station** (1998) ★ 8.0 /10 31,520 Rate This

Central do Brasil (*original title*)
R | 1h 53min | Drama | 20 November 1998 (USA)

CENTRAL STATION
Movie poster showing a woman and a child in a doorway.

1:54 | Trailer 1 VIDEO | 22 IMAGES

a On Disc at Amazon

An emotive journey of a former school teacher, who writes letters for illiterate people, and a young boy, whose mother has just died, as they search for the father he never knew.

Director: Walter Salles
Writers: Marcos Bernstein, João Emanuel Carneiro | 1 more credit »
Stars: Fernanda Montenegro, Vinícius de Oliveira, Marília Pêra | See full cast & crew »

Metascore 80 From metacritic.com | Reviews | Prime Video Watch Now From \$2.99 (SD) on Prime Video



FULL CAST AND CREW | TRIVIA | USER REVIEWS | IMDbPro | MORE ▾ SHARE

+ **Star Wars: The Last Jedi** (2017) ★ 7.3 /10 404,499 Rate This

Star Wars: Episode VIII - The Last Jedi (*original title*)
PG-13 | 2h 32min | Action, Adventure, Fantasy | 15 December 2017 (USA)

STAR WARS: THE LAST JEDI
Movie poster showing Rey, Luke Skywalker, and other characters.

Rey develops her newly discovered abilities with the guidance of Luke Skywalker, who is unsettled by the strength of her powers. Meanwhile, the Resistance prepares for battle with the First Order.

Director: Rian Johnson
Writers: Rian Johnson, George Lucas (based on characters created by)
Stars: Daisy Ridley, John Boyega, Mark Hamill | See full cast & crew »

Metascore 85 From metacritic.com | Reviews | Popularity 84 (▲ 3)

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

ON DISC

Same pred may corr. to diff DOM tree nodes

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

- Challenge: slight variations from page to page

Central Station (1998)
Central do Brasil (original title)
R | 1h 53min | Drama | 20 November 1998 (USA)

The Fog of War: Eleven Lessons from the Life of Robert S. McNamara (2003)
PG-13 | 1h 47min | Documentary, Biography, History | 5 March 2004 (USA)

On Disc at Amazon

An emotive journey of a former school teacher, who writes letters for illiterate people, and a young boy, whose mother has just died, as they search for the father he never knew.

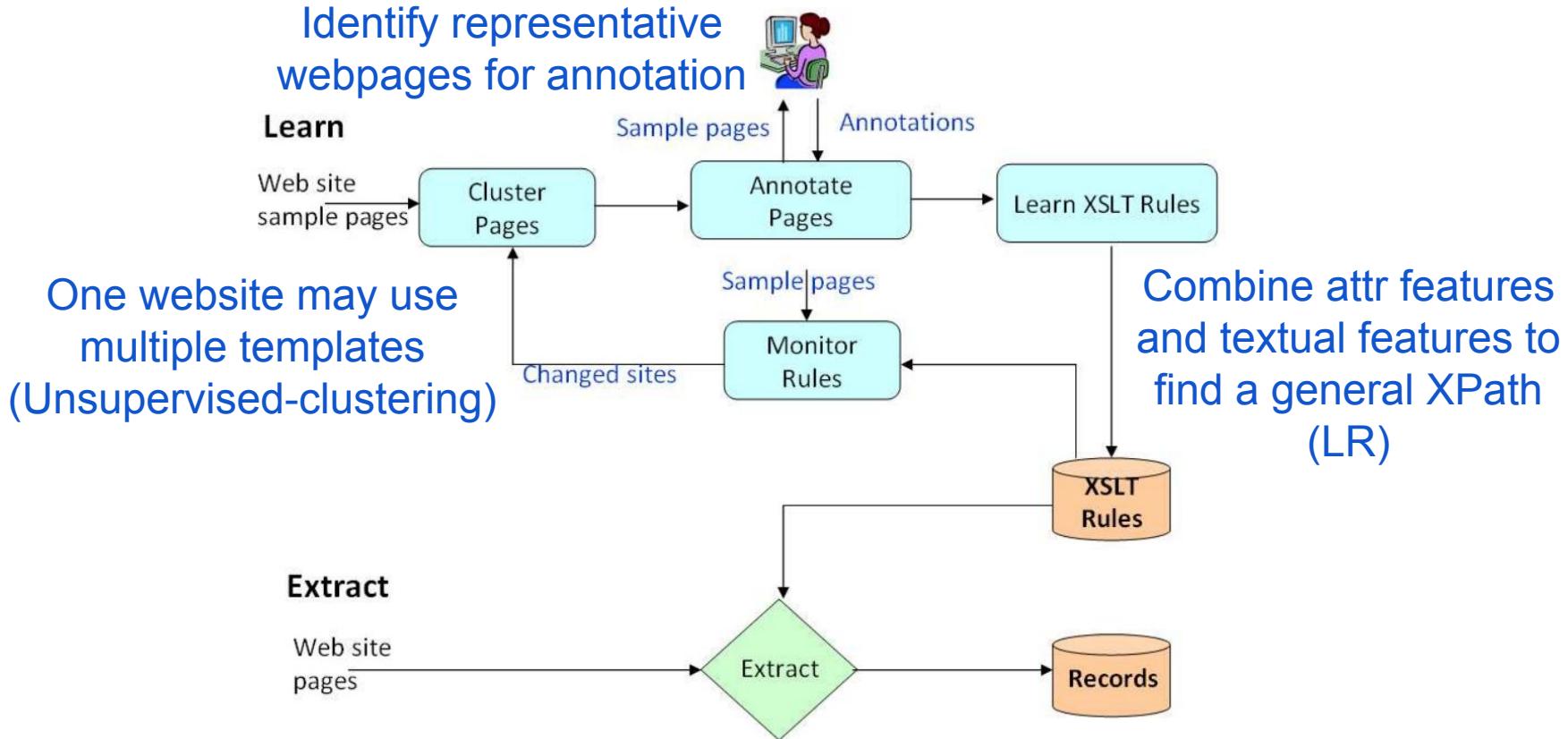
Director: [Walter Salles](#)
Writers: [Marcos Bernstein](#), [João Emanuel Carneiro](#) | 1 more credit »
Stars: [Fernanda Montenegro](#), [Vinícius de Oliveira](#), [Marília Pêra](#) | See full cast & crew »

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

ON DISC

Same DOM tree node may correspond to diff preds

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



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

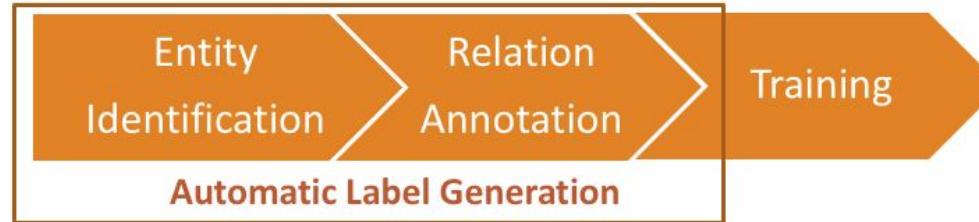
Ensure high precision

Distantly Supervised Extraction

2013 (Deep ML)

- **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

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



Movie entity



Runtime



Genre Release Date

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")

Talk on Tue Research 11am, Poster on Tue Posters | 17:30pm

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

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

Very high precision

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]

2013 (Deep ML)



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

Deep learning

- Use RNN, CNN, attention for RE
- Data programming / Heterogeneous learning
- Revisit DOM extraction

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 production credits and technical specifications. The production credits section includes 'Company Credits' and a list of production companies: Lucasfilm, Walt Disney Pictures, Allison Shearmur Productions, with a 'See more' link. Below it is a 'Show more on IMDbPro' link. The technical specs section lists runtime (135 min), sound mix options (Dolby Atmos, DTS, DTS:X, 12-Track Digital Sound, Auro 11.1, Dolby Digital, Dolby Surround 7.1), color (Color), and aspect ratio (2.39 : 1). A 'See full technical specs' link 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 »

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

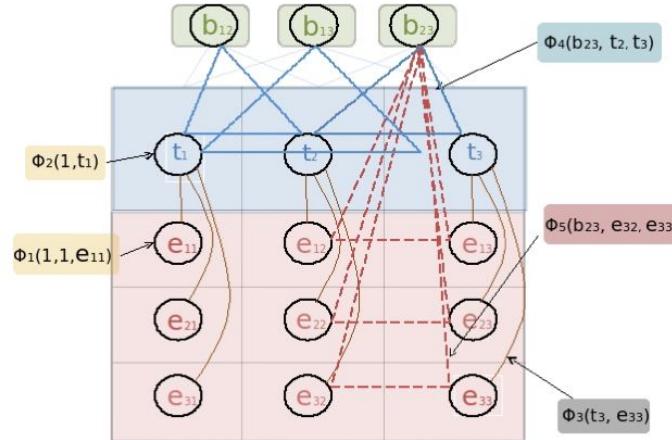
- Model table annotation using interrelated random variables, represented by a probabilistic graphical model
 - Cell text (in Web table) and entity label (in catalog)
 - Column header (in Web table) and type label (in catalog)
 - Column type and cell entity (in Web table)

Extraction from semi-structured data

- WebTables: search, extraction
- DOM tree: wrapper induction

2008 (Semi-stru)

Check-out 10-Year Best Paper Award for WebTable Search on Thursday!



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

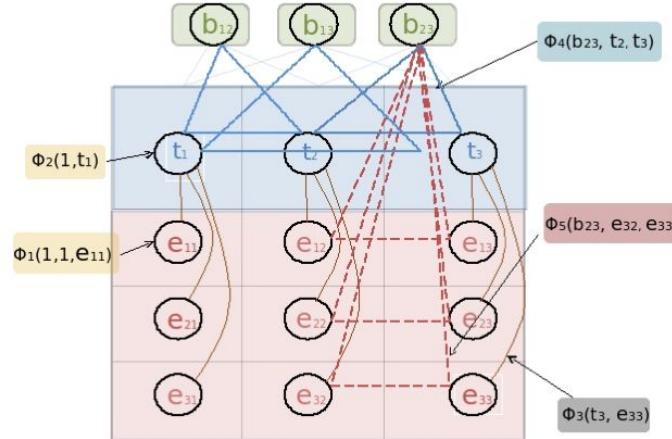
- Model table annotation using interrelated random variables, represented by a probabilistic graphical model
 - Pair of column types (in Web table) and relation (in catalog)
 - Entity pairs (in Web table) and relation (in catalog)

Extraction from semi-structured data

- WebTables: search, extraction
- DOM tree: wrapper induction

2008 (Semi-stru)

Check-out 10-Year Best Paper Award for WebTable Search on Thursday!

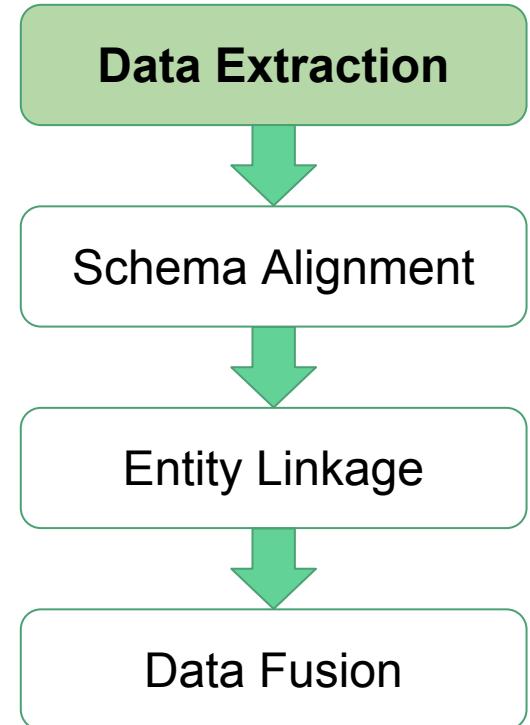


Challenges in Applying ML on DX

- 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?

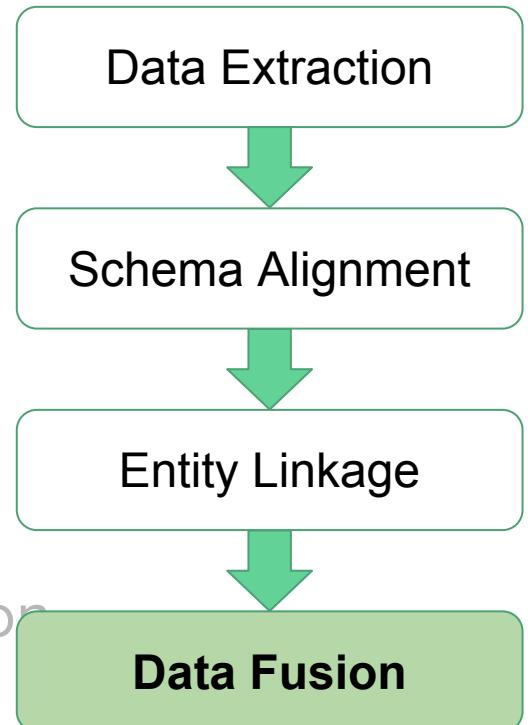
Recipe for Data Extraction

- 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**
 - **DL effective for texts and LR is often effective for semi-stru data**



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 (2,020 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 rivers in the United States								
#	Name	Mouth ^[5]	Length	Source coordinates ^[11]	Mouth coordinates ^[11]	Watershed area ^[12]	Discharge ^[12]	States, provinces, and image ^{[5][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 ⁵ , North Dakota, South Dakota, Nebraska, Iowa, Kansas, Missouri ^[16] 
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 ⁵ , Wisconsin, Iowa, Illinois, Missouri, Kentucky, Tennessee, Arkansas, Mississippi, Louisiana ^[16] 

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

Source reports
a value for a fact

Conflicting value

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

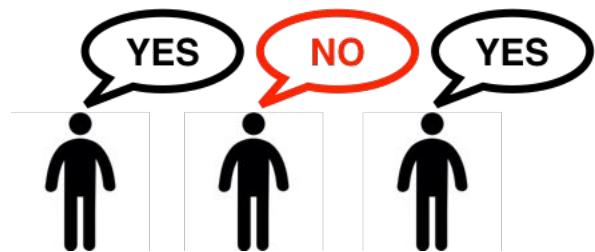
Source	River	Attribute	Value
KG	Mississippi River	Length	2,320 mi
KG	Missouri River	Length	2,341 mi
Wikipedia	Mississippi River	Length	2,202 mi
Wikipedia	Missouri River	Length	2,341 mi
USGS	Mississippi River	Length	2,340 mi
USGS	Missouri River	Length	2,540 mi

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

~1996 (Rule-based)

1979

(Statistical learning)

2016 (Deep ML)

2007 (Probabilistic)

Domain-specific Strategies

- Keep all values
- Pick a random value
- Take the average value
- Take the most recent value
- ...

Deep learning

- Use Restricted Boltzmann Machine; one layer version is equivalent with Dawid-Skene model
- Knowledge graph embeddings

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.

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

Normalizing constant

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

error-rate scores
(model parameters)

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.

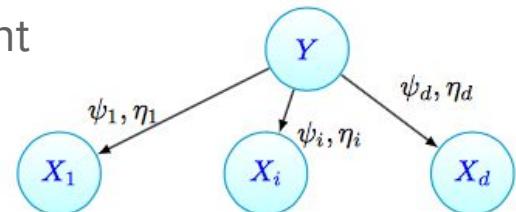
But the number of sources can be in the thousands or millions and training data is limited!

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.

An Intro in Probabilistic Graphical Models

Bayesian Networks (BNs)

Local Markov Assumption: A variable X is independent of its non-descendants given its parents (and *only* its parents).

An Intro in Probabilistic Graphical Models

Bayesian Networks (BNs)

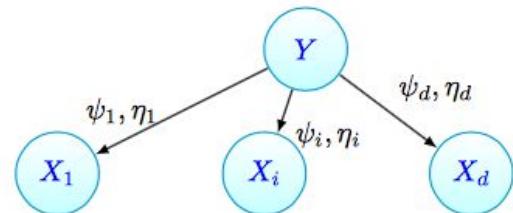
Local Markov Assumption: A variable X is independent of its non-descendants given its parents (and *only* its parents).

Recipe for BNs

Set of random variables X

Directed acyclic graph (each $X[i]$ is a vertex)

Conditional probability tables $P(X \mid \text{Parents}(X))$



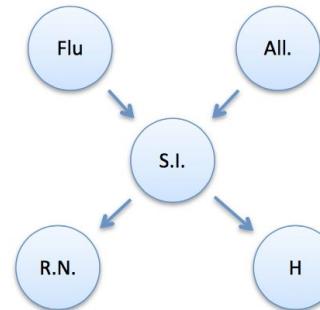
Joint distribution: Factorizes over conditional probability tables

An Intro in Probabilistic Graphical Models

Where do independence assumptions come from?

Causal structure captures domain knowledge

- The flu causes sinus inflammation
- Allergies *also* cause sinus inflammation
- Sinus inflammation causes a runny nose
- Sinus inflammation causes headaches

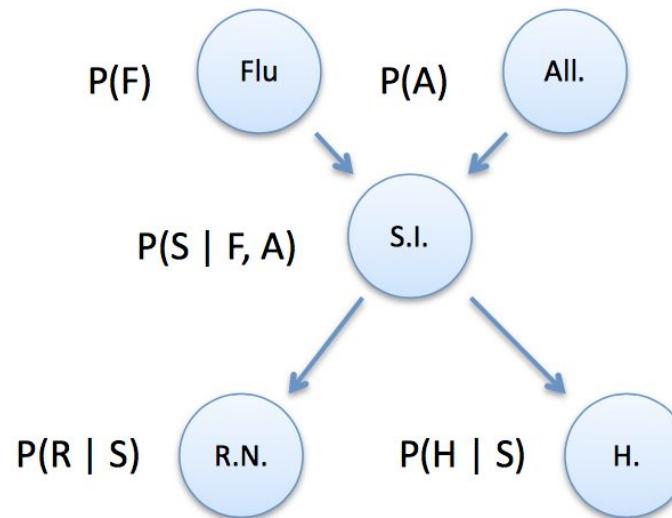


[Example by Andrew McCallum]

An Intro in Probabilistic Graphical Models

Factored joint distribution

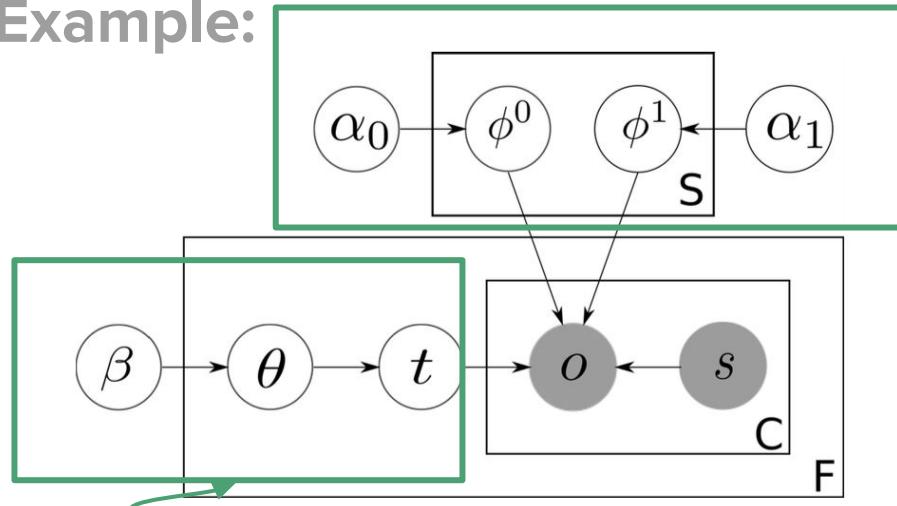
$$\begin{aligned} & P(F, A, S, R, H) \\ & = P(F) \\ & \quad P(A) \\ & \quad P(S | F, A) \\ & \quad P(R | S) \\ & \quad P(H | S) \end{aligned}$$



[Example by Andrew McCallum]

Probabilistic Graphical Models for Data Fusion

Example:



Source
Quality

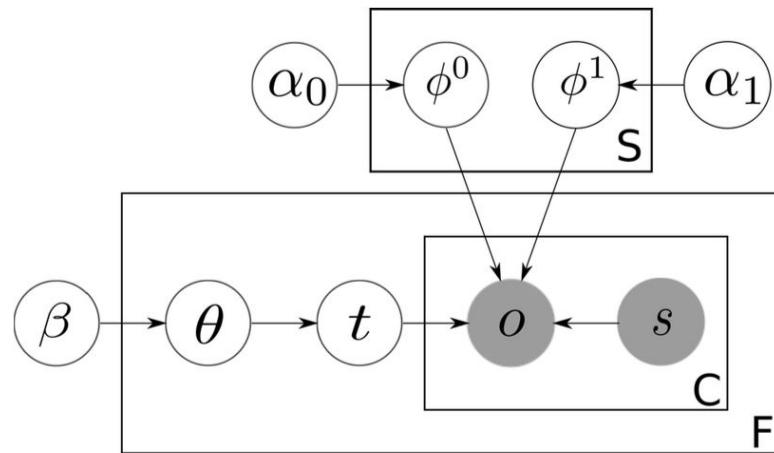
**Setup: Identify true
source claims**

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
...

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

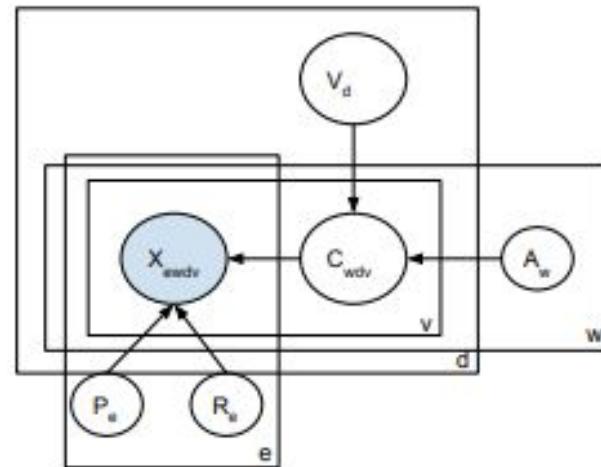
Extensive work on modeling source observations and source interactions to address limitations of basic Dawid-Skene.

Probabilistic Graphical Models for Data Fusion



[Zhao et al., VLDB 2012]

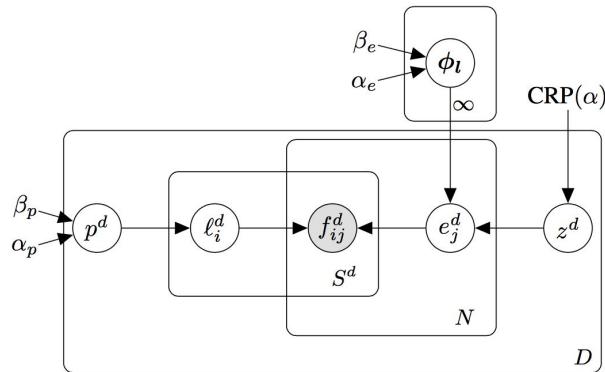
Modeling both source quality
and extractor accuracy



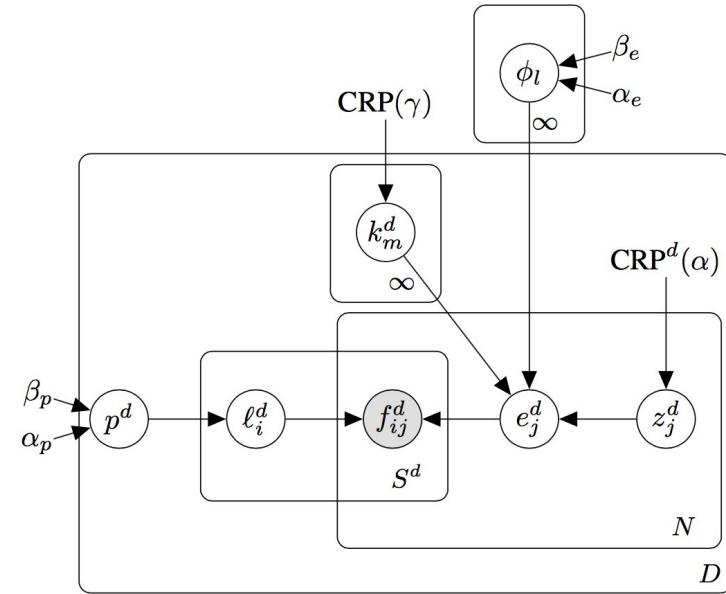
[Dong et al., VLDB 2015]

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		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.

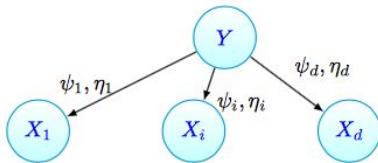
PGMs in Data Fusion [Li et al., VLDB'14]

Category	Method	Stock				Flight			
		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
	ACCU SIM	.918	.913	.17	-.16	.903	.844	.2	-.09
	ACCU FORMAT	.918	.911	.17	-.16	.903	.844	.2	-.09
	ACCU SIM ATTR	.950	.929	.17	-.16	.952	.833	.19	-.08
	ACCU FORMAT ATTR	.948	.930	.17	-.16	.952	.833	.19	-.08
Copying affected	ACCU COPY	.958	.892	.28	-.11	.960	.943	.16	-.14

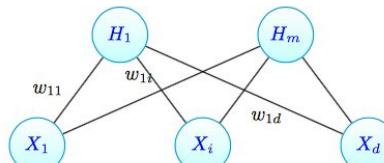
Modeling the quality of data sources leads to improved accuracy.

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

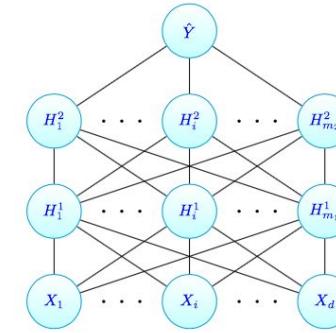
Theorem: The Dawid and Skene model is *equivalent* to a Restricted Boltzmann Machine (RBM) with a single hidden node.



Dawid and Skene model.



A RBM with d visible and m hidden units.

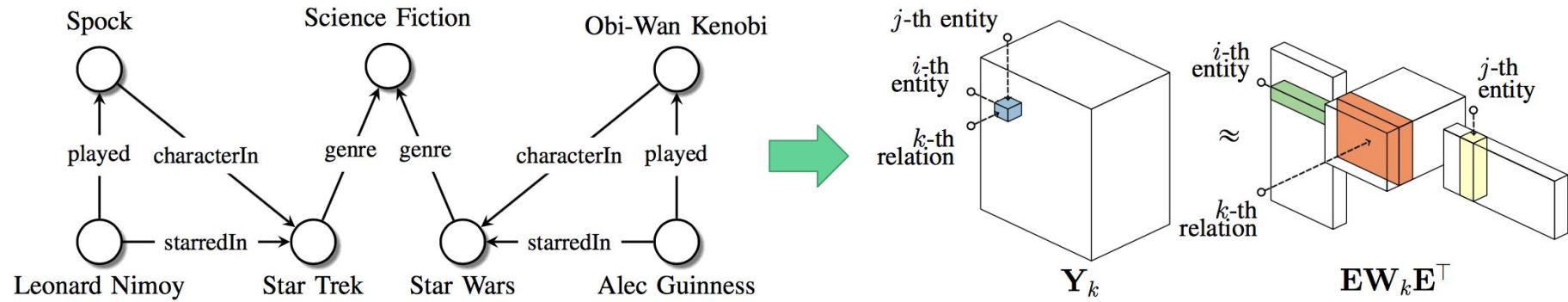


Sketch of a two-hidden-layer RBM-based DNN.

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

Knowledge Graph Embeddings

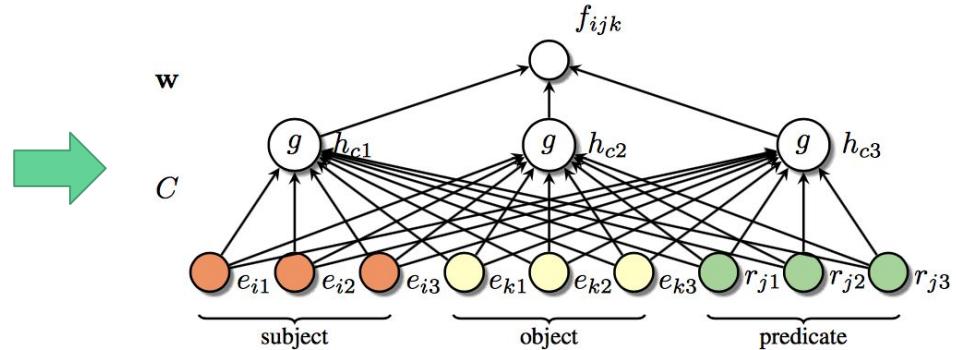
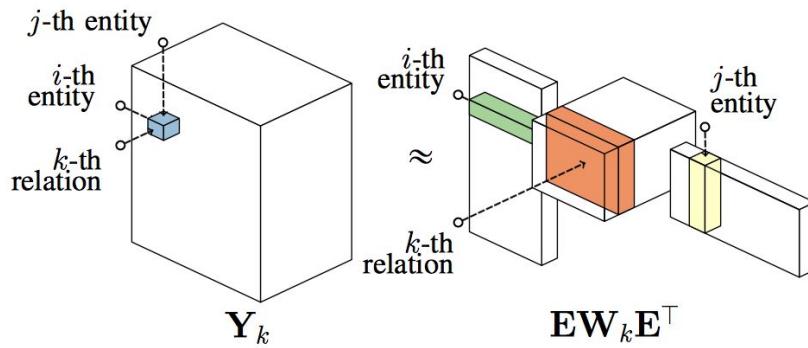
[Survey: Nicket et al., 2015]



A knowledge graph can be encoded as a tensor.

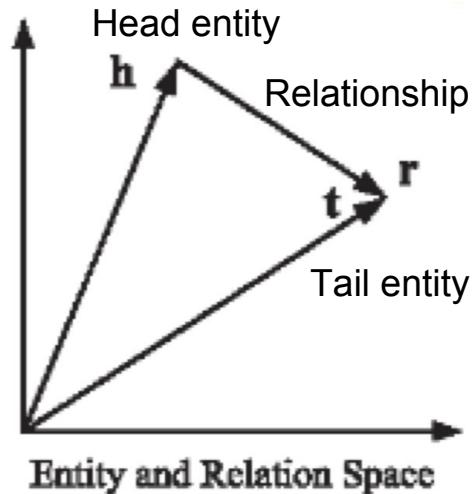
Knowledge Graph Embeddings

[Survey: Nicket et al., 2015]



Neural networks can be used to obtain richer representations.

Knowledge Graph Embeddings



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

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
...

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]

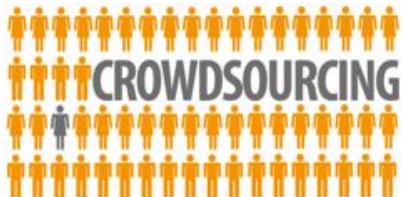
Limit the informative parameters of the model by using domain knowledge

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

Example:

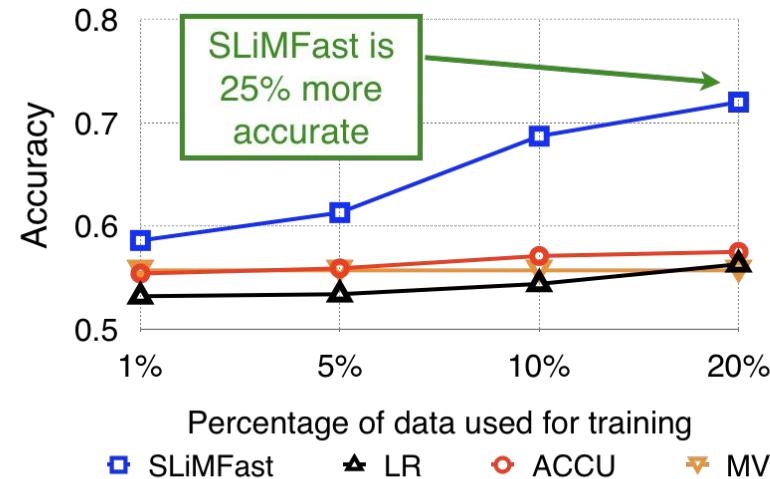
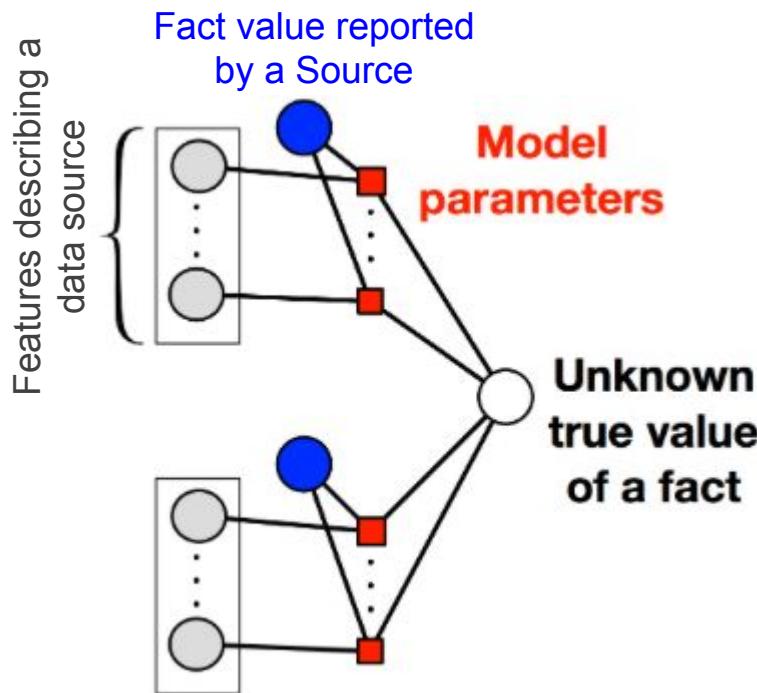


- newly registered similar to existing domain
 - traffic statistics
 - text quality (e.g., misspelled words, grammatical errors)
 - sentiment analysis
-
- avg. time per task
 - number of tasks
 - market used



SLIMFast: Discriminative Data Fusion

[Rekatsinas et al., SIGMOD'17]



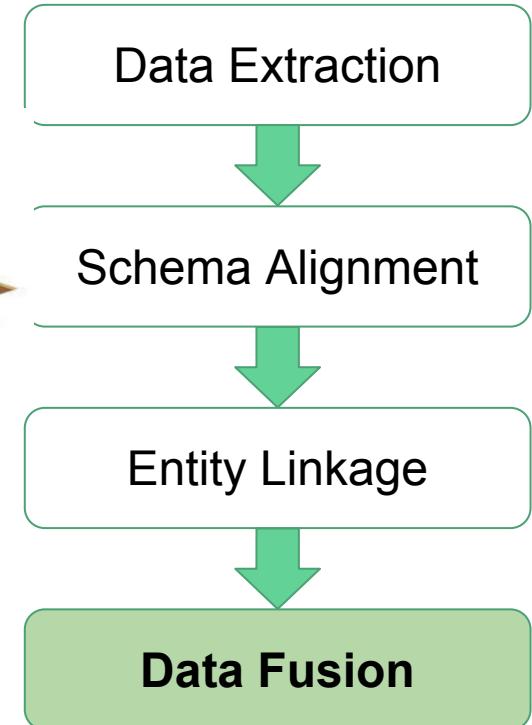
Genomics data: 2.7k sources (articles), 571 objects (gene-disease), 4 domain features (year, citation, author, journal)

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?

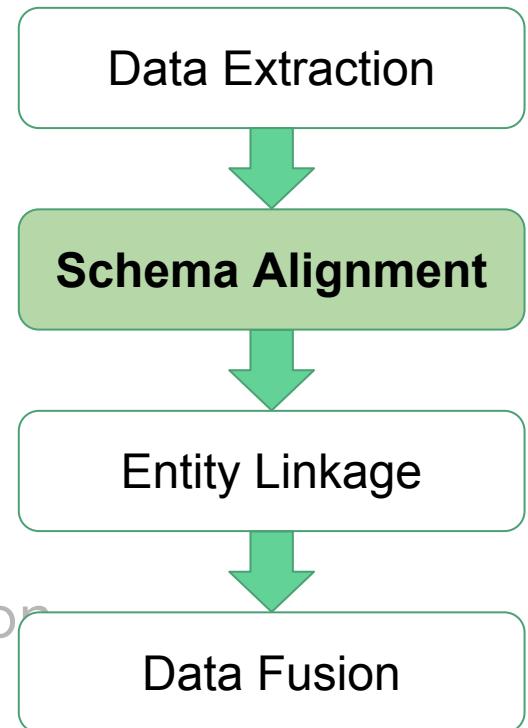
Recipe for Data Fusion

- 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
- Part IV. Conclusions and research direction



What is Schema Alignment?

- Definition: Align schemas and understand which attributes have the same semantics.

IMDB



Anahí
Actress | Music Department | Soundtrack

SEE RANK

Anahí 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 (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	<small>edit</small>
▼ 1 reference	<small>imported from</small>	Italian Wikipedia
	<small>+ add reference</small>	
	<small>+ add value</small>	

Quick Tour for Schema Alignment

S1	(name, hPhone, hAddr, oPhone, oAddr)
S2	(name, phone, addr, email)
S3	a: (id, name); b: (id, resPh, workPh)
S4	(name, pPh, pAddr)
S5	(name, wPh, wAddr)

Mediated Schema

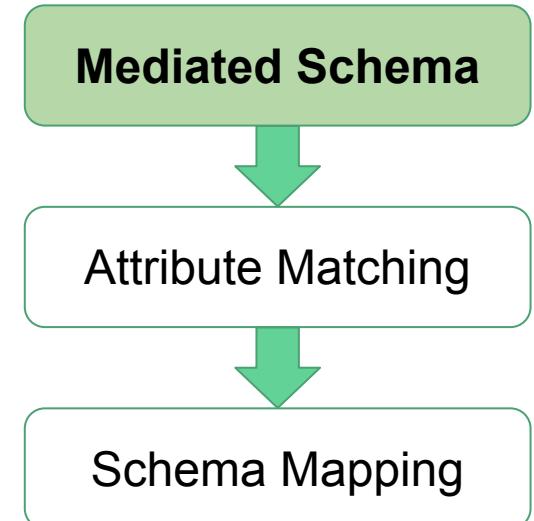
Attribute Matching

Schema Mapping

Quick Tour for Schema Alignment

- **Mediated schema:** a unified and virtual view of the salient aspects of the domain

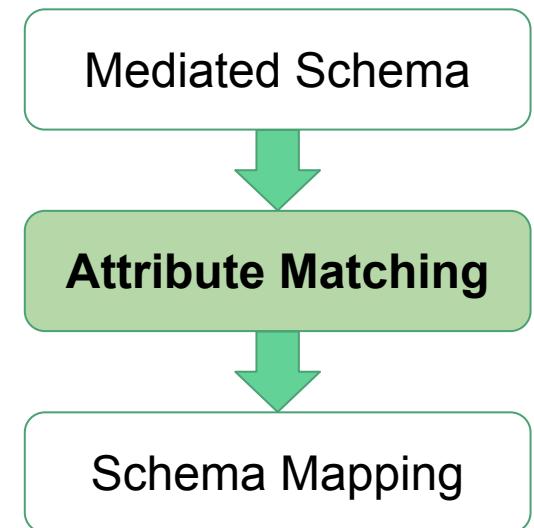
S1	(name, hPhone, hAddr, oPhone, oAddr)
S2	(name, phone, addr, email)
S3	a: (id, name); b: (id, resPh, workPh)
S4	(name, pPh, pAddr)
S5	(name, wPh, wAddr)
MS	(n, pP, pA, wP, wA)



Quick Tour for Schema Alignment

- **Attribute matching:** correspondences between schema attributes

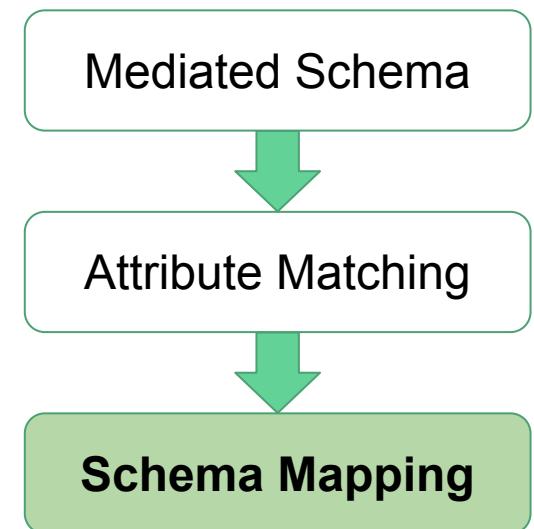
S1	(name, hPhone, hAddr, oPhone, oAddr)
S2	(name, phone, addr, email)
S3	a: (id, name); b: (id, resPh, workPh)
S4	(name, pPh, pAddr)
S5	(name, wPh, wAddr)
MS	(n, pP, pA, wP, wA)
MSAM	MS.n: S1.name, S2.name, S3a.name, ... MS.pP: S1.hPhone, S3b.resPh, S4.pPh MS.pA: S1.hAddr, S4.pAddr MS.wP: S1.oPhone, S2.phone, ... MS.wA: S1.oAddr, S2.addr, S5.wAddr



Quick Tour for Schema Alignment

- **Schema mapping:** transformation between records in different schemas

S1	(name, hPhone, hAddr, oPhone, oAddr)
S2	(name, phone, addr, email)
S3	a: (id, name); b: (id, resPh, workPh)
S4	(name, pPh, pAddr)
S5	(name, wPh, wAddr)
MS	(n, pP, pA, wP, wA)
MSSM (GAV)	MS(n, pP, pA, wP, wA) :- S1(n, pP, pA, wP, wA) MS(n, _, _, wP, wA) :- S2(n, wP, wA, e) MS(n, pP, _, wP, _) :- S3a(i, n), S3b(i, pP, wP) MS(n, pP, pA, _, _) :- S4(n, pP, pA) MS(n, _, _, wP, wA) :- S5(n, wP, wA)



30 Years of Schema Alignment

Description Logics

- Gav vs. Lav. vs. Glav
- Answering queries using views
- Warehouse vs. EII



~1990 (Desc Logics)

1994 (Early ML)

2005 (Dataspaces)

2013 (Deep ML)

Semi-Auto mapping

- Learning to match
- Schema mapping: Clio
- Data exchange

Logic & Deep learning

- Collective disc. by PSL
- Universal schema

Early ML Models

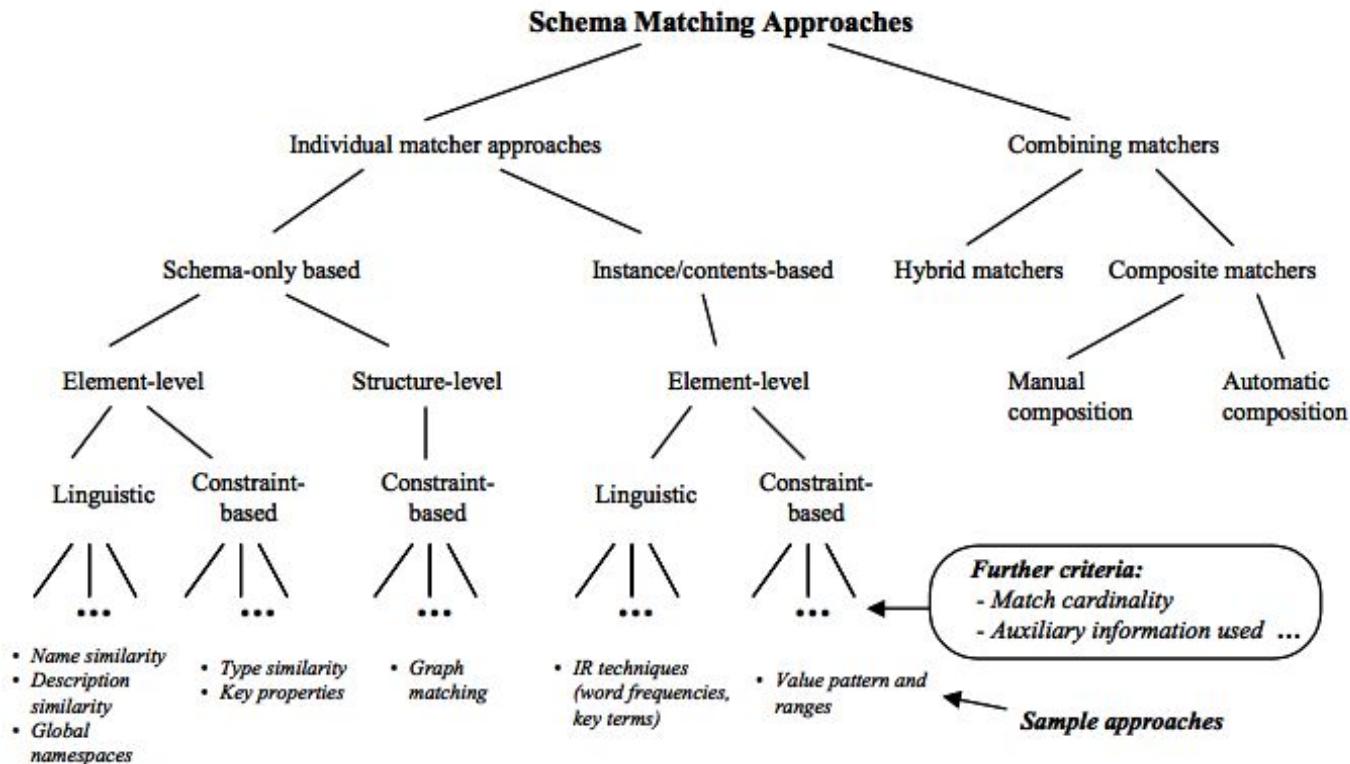
[Rahm and Bernstein, VLDBJ'2001]

~2000 (Early ML)



Semi-Auto mapping

- Learning to match
- Schema mapping: Clio
- Data exchange



Further criteria:
- Match cardinality
- Auxiliary information used ...

Sample approaches

Signals: name, description, type, key, graph structure, values

Early ML Models

[Doan et al., Sigmod'01]

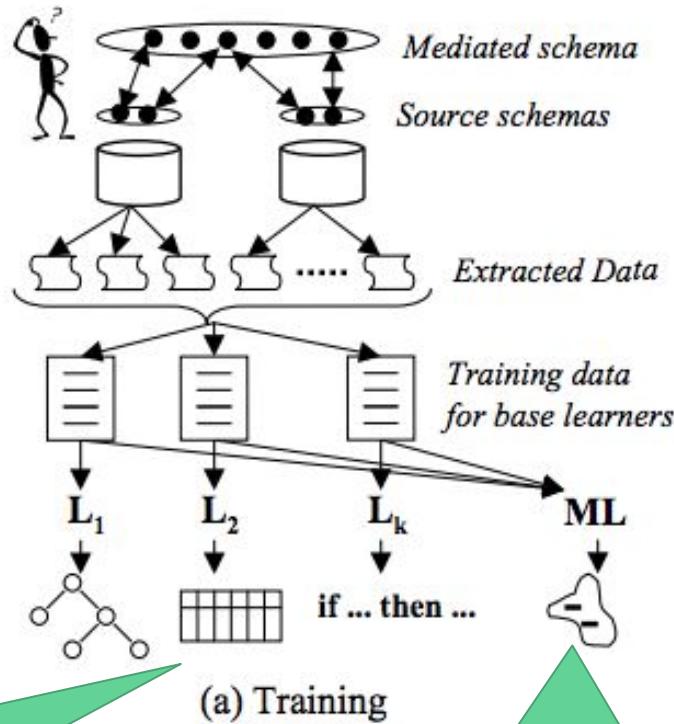
~2000 (Early ML)



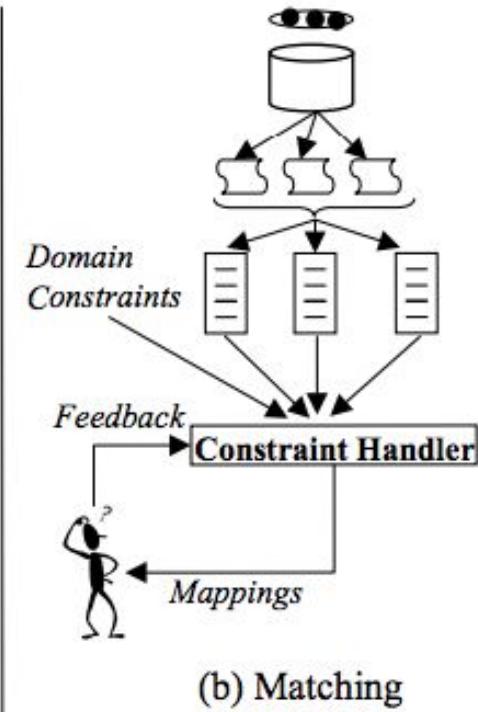
Semi-Auto mapping

- Learning to match
- Schema mapping: Clio
- Data exchange

Base learners: kNN, naive Bayes, etc.



Meta learner--Stacking



Early ML Models

[Doan et al., Sigmod'01]

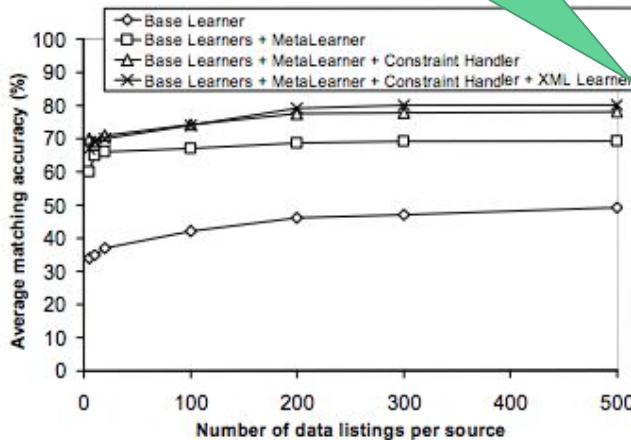
~2000 (Early ML)



Semi-Auto mapping

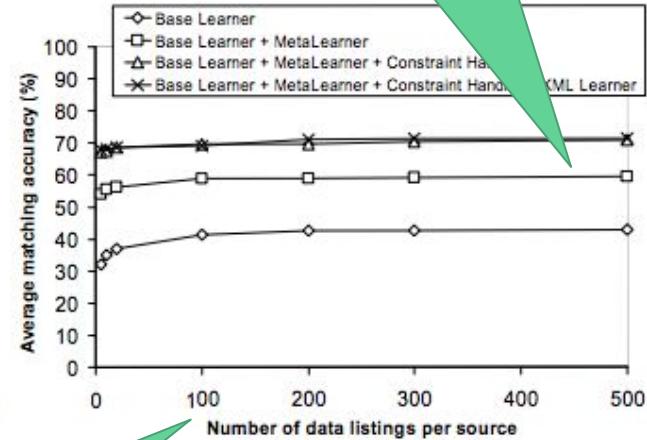
- Learning to match
- Schema mapping: Clio
- Data exchange

Avg Accuracy: 71-92%



(b) Matching accuracy for Real Estate I

Meta learning and constraints help



(c) Matching accuracy for Time Schedule

More data instances help

Collective Mapping Discovery by PSL [Kimmig et al, ICDE'17]

Step 1. Generate candidate mappings

E.g., $\theta_0 : \text{proj}(t, m, l) \wedge \text{emp}(m, n, c) \rightarrow \exists o. \text{task}(t, n, o)$

$\theta_1 : \text{proj}(t, m, l) \wedge \text{emp}(l, n, c) \rightarrow \exists o. \text{task}(t, n, o)$

$\theta_2 : \text{proj}(t, m, l) \wedge \text{emp}(m, n, c) \rightarrow \exists o. \text{task}(t, n, o) \wedge \text{org}(o, c)$

$\theta_3 : \text{proj}(t, m, l) \wedge \text{emp}(l, n, c) \rightarrow \exists o. \text{task}(t, n, o) \wedge \text{org}(o, c)$

2013 (Deep ML)



Logic & Deep learning

- Collective disc. by PSL
- Universal schema

Step 2. Solve PSL

1. Prefer fewer mappings: penalty = #atoms

$\text{size}_m(F) : in(F) \rightarrow \perp$

$1 : J(T) \rightarrow \exists F. \text{covers}(F, T) \wedge in(F)$

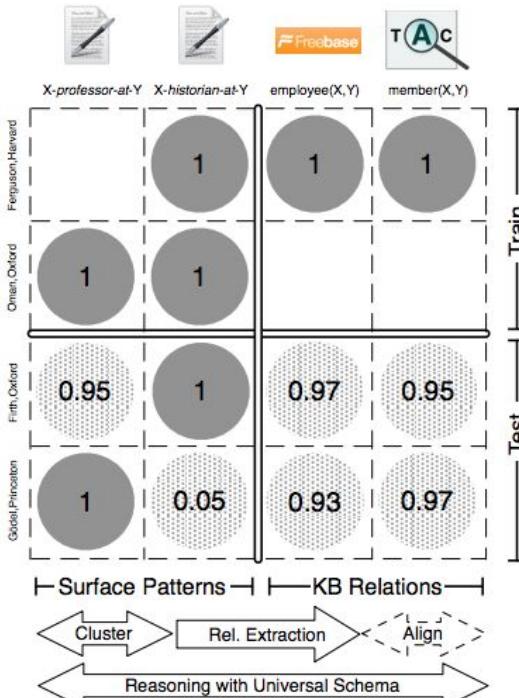
$1 : in(F) \wedge \text{creates}(F, T) \rightarrow J(T)$

3. Tuples inferred from the mapping should exist

2. An existing tuple can be inferred from the mappings

Universal Schema [Riedel et al., NAACL'13][Yao et al., AKBC'13]

- Attribute matching → Instance inference



2013 (Deep ML)

Logic & Deep learning

- Collective disc. by PSL
- Universal schema

Matrix factorization

	per/actor	loc/country	lawyer	...	company
Barack Obama		1			
Ruth B. Ginsburg		1			
New York					
Argentina	.89				
Brad Pitt	1				
IBM					1
...					

Type prediction

Relation prediction

Universal Schema [Riedel et al., NAACL'13]

- Attribute matching → Instance inference
- $f(e_s, r, e_o)$ is computed using embeddings; the higher, the more likely to be true
- DistMult is a relation embedding model

2013 (Deep ML)



Logic & Deep learning

- Collective disc. by PSL
- Universal schema

Limitation: Cannot apply to new entities or relations

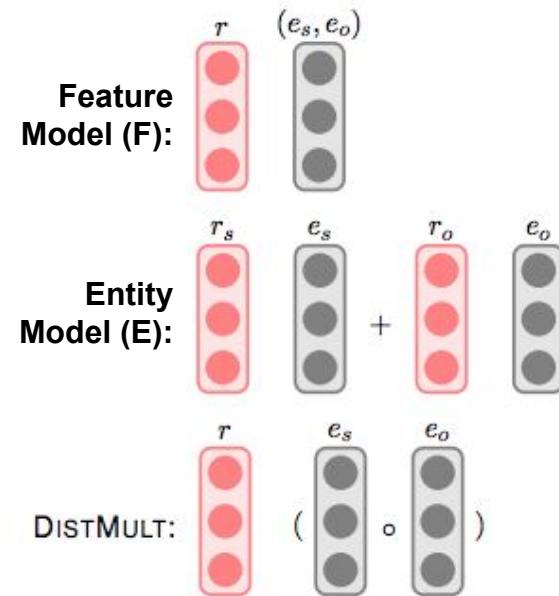


Figure 3: The continuous representations for model F, E and DISTMULT. [Toutanova et al., EMNLP'15]

Columnless Univ. Schema w. CNN [Toutanova et al., EMNLP'15]

- Relation: organizationFoundedBy

Textual Pattern	Count
SUBJECT $\xrightarrow{\text{appos}}$ founder $\xrightarrow{\text{prep}}$ of $\xrightarrow{\text{pobj}}$ OBJECT	12
SUBJECT $\xleftarrow{\text{nsubj}}$ co-founded $\xrightarrow{\text{dobj}}$ OBJECT	3
SUBJECT $\xrightarrow{\text{appos}}$ co-founder $\xrightarrow{\text{prep}}$ of $\xrightarrow{\text{pobj}}$ OBJECT	
SUBJECT $\xrightarrow{\text{conj}}$ co-founder $\xrightarrow{\text{prep}}$ of $\xrightarrow{\text{pobj}}$ OBJECT	
SUBJECT $\xleftarrow{\text{pobj}}$ with $\xrightarrow{\text{prep}}$ co-founded $\xrightarrow{\text{dobj}}$ OBJECT	
SUBJECT $\xleftarrow{\text{nsubj}}$ signed $\xrightarrow{\text{xcomp}}$ establishing $\xrightarrow{\text{dobj}}$ OBJECT	
SUBJECT $\xleftarrow{\text{pobj}}$ with $\xrightarrow{\text{prep}}$ founders $\xrightarrow{\text{prep}}$ of $\xrightarrow{\text{pobj}}$ OBJECT	2
SUBJECT $\xrightarrow{\text{appos}}$ founders $\xrightarrow{\text{prep}}$ of $\xrightarrow{\text{pobj}}$ OBJECT	2
SUBJECT $\xleftarrow{\text{nsubj}}$ one $\xrightarrow{\text{prep}}$ of $\xrightarrow{\text{pobj}}$ founders $\xrightarrow{\text{prep}}$ of $\xrightarrow{\text{pobj}}$ OBJECT	2
SUBJECT $\xleftarrow{\text{nsubj}}$ founded $\xrightarrow{\text{dobj}}$, production $\xrightarrow{\text{conj}}$ OBJECT	2
SUBJECT $\xleftarrow{\text{appos}}$ partner $\xrightarrow{\text{pobj}}$ with $\xrightarrow{\text{prep}}$ founded $\xrightarrow{\text{dobj}}$ production $\xrightarrow{\text{conj}}$ OBJECT	2
SUBJECT $\xleftarrow{\text{pobj}}$ by $\xrightarrow{\text{prep}}$ co-founded $\xrightarrow{\text{rcmod}}$ OBJECT	1
SUBJECT $\xleftarrow{\text{nn}}$ co-founder $\xrightarrow{\text{prep}}$ of $\xrightarrow{\text{pobj}}$ OBJECT	1
SUBJECT $\xrightarrow{\text{dep}}$ co-founder $\xrightarrow{\text{prep}}$ of $\xrightarrow{\text{pobj}}$ OBJECT	1
SUBJECT $\xleftarrow{\text{nsubj}}$ helped $\xrightarrow{\text{xcomp}}$ establish $\xrightarrow{\text{dobj}}$ OBJECT	1
SUBJECT $\xleftarrow{\text{nsubj}}$ signed $\xrightarrow{\text{xcomp}}$ creating $\xrightarrow{\text{dobj}}$ OBJECT	1

Similarity of phrases
→ CNN

2013 (Deep ML)



Logic & Deep learning

- Collective disc. by PSL
- Universal schema

Columnless Univ. Schema w. CNN

[Toutanova et al., EMNLP'15]

2013 (Deep ML)



Logic & Deep learning

- Collective disc. by PSL
- Universal schema

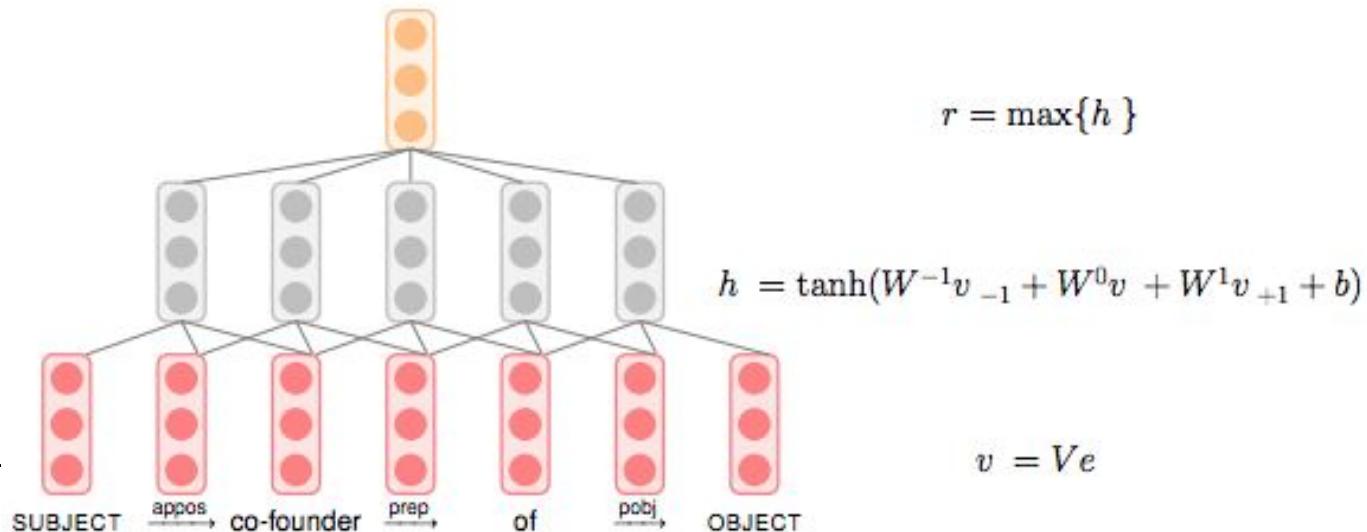


Figure 4: The convolutional neural network architecture for representing textual relations.

Columnless Univ. Schema w. RNN [Verga et al., ACL'16]

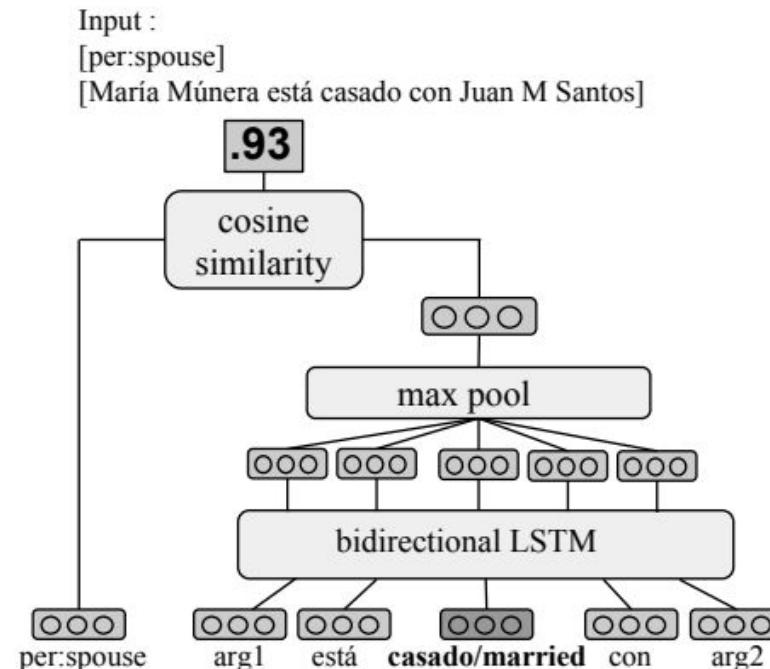
- Similar sequences of context tokens should be embedded similarly

2013 (Deep ML)



Logic & Deep learning

- Collective disc. by PSL
- Universal schema



Rowless Univ. Schema [Verga et al., ACL'16]

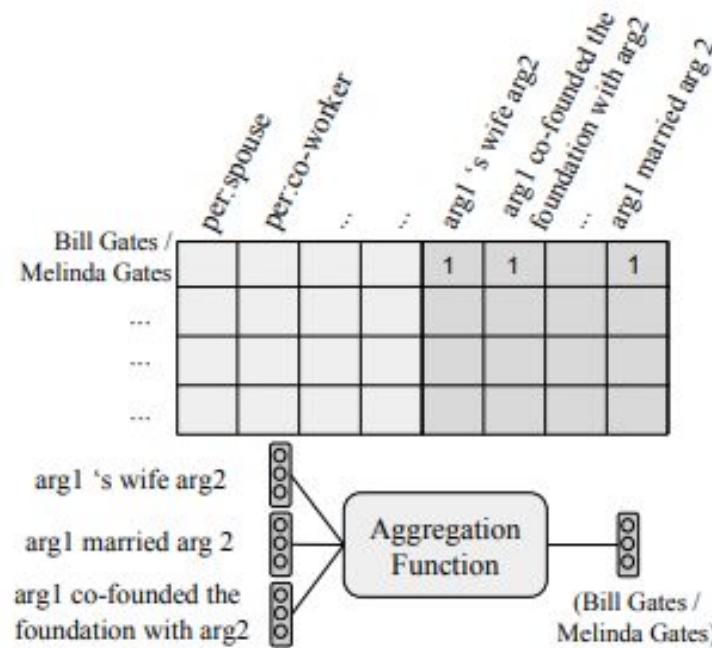
- Infer relation from a set of observed relations
- Similar to schema mapping w. signals from values

2013 (Deep ML)



Logic & Deep learning

- Collective disc. by PSL
- Universal schema



Rowless Univ. Schema [Verga et al., ACL'16]

2013 (Deep ML)

Rowless & Columnless

Model	MRR	Hits@10
Entity-pair Embeddings	31.85	51.72
Entity-pair Embeddings-LSTM	33.37	54.39
Attention	31.92	51.67
Attention-LSTM	30.00	53.35
Max Relation	31.71	51.94
Max Relation-LSTM	30.77	54.80

(a)

Model	MRR	Hits@10
Entity-pair Embeddings	5.23	11.94
Attention	29.75	49.69
Attention-LSTM	27.95	51.05
Max Relation	28.46	48.15
Max Relation-LSTM	29.61	54.19

(b)

Logic & Deep learning

- Collective disc. by PSL
- Universal schema

Recall still low

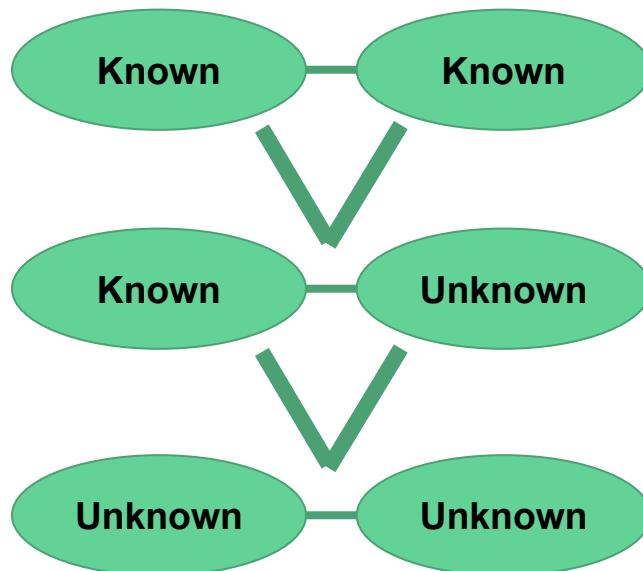
Similar for new entity pairs

Schema Mapping vs. Universal Schema

	Schema matching	Universal schema
<i>Granularity</i>	Column-level decision	Cell-level decision
<i>Expressiveness</i>	Mainly 1:1 mapping	Allow overlap, subset/superset, etc.
<i>Signals</i>	Name, description, type, key, graph structure, values	Values
<i>Results</i>	Accu: 70-90%	MRR=~0.3, Hits@10=~0.5
<i>Community</i>	Database	NLP

Challenges in Applying Deep Learning on SM

- How can we combine techs from schema matching and universal schema??



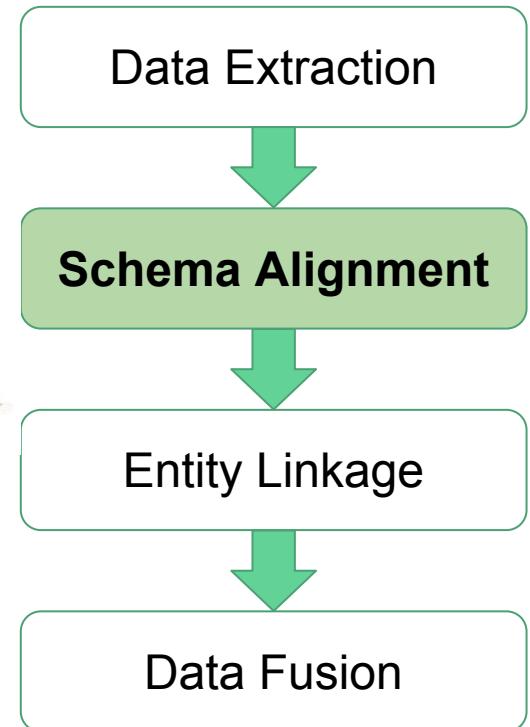
Leverage knowledge by inference

Leverage knowledge on types

Rowless

Recipe for Schema Alignment

- Problem definition: **Align attributes with the same semantics**
- Short answers
 - **Interactive semi-automatic mapping**
 - **DL-based universal schema revived the field**
 - **Combine schema matching and universal schema for future**

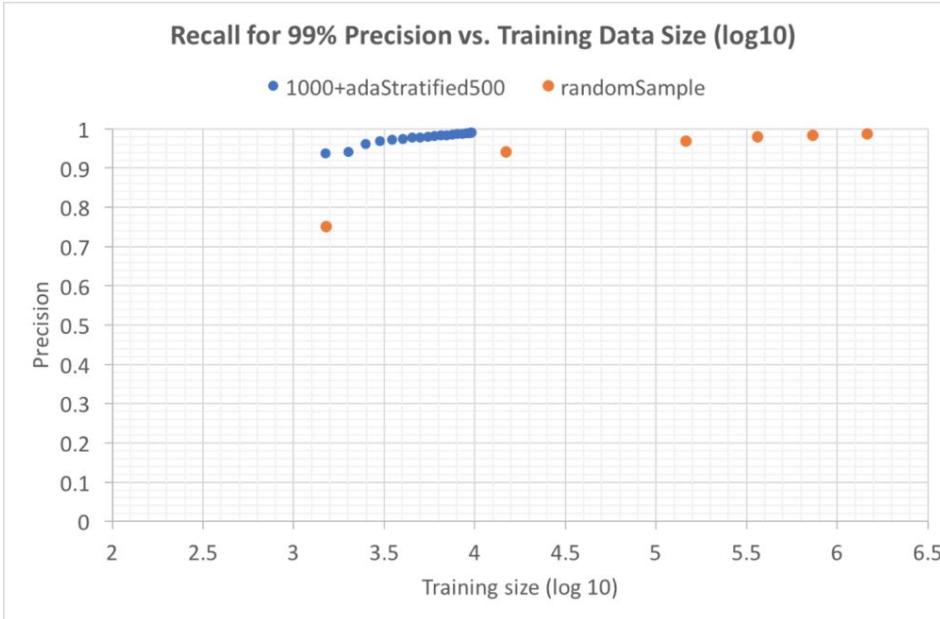


Revisit Theme I. Which ML Model Works Best?

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				X	
Text extraction	X	X		X		X
Schema alignment	X		X	X	X	X

For structured data, RF works well, and LR is often effective
For texts and semantics, deep learning shows big promise

Revisit Theme II. Does Supervised Learning Apply to DI?

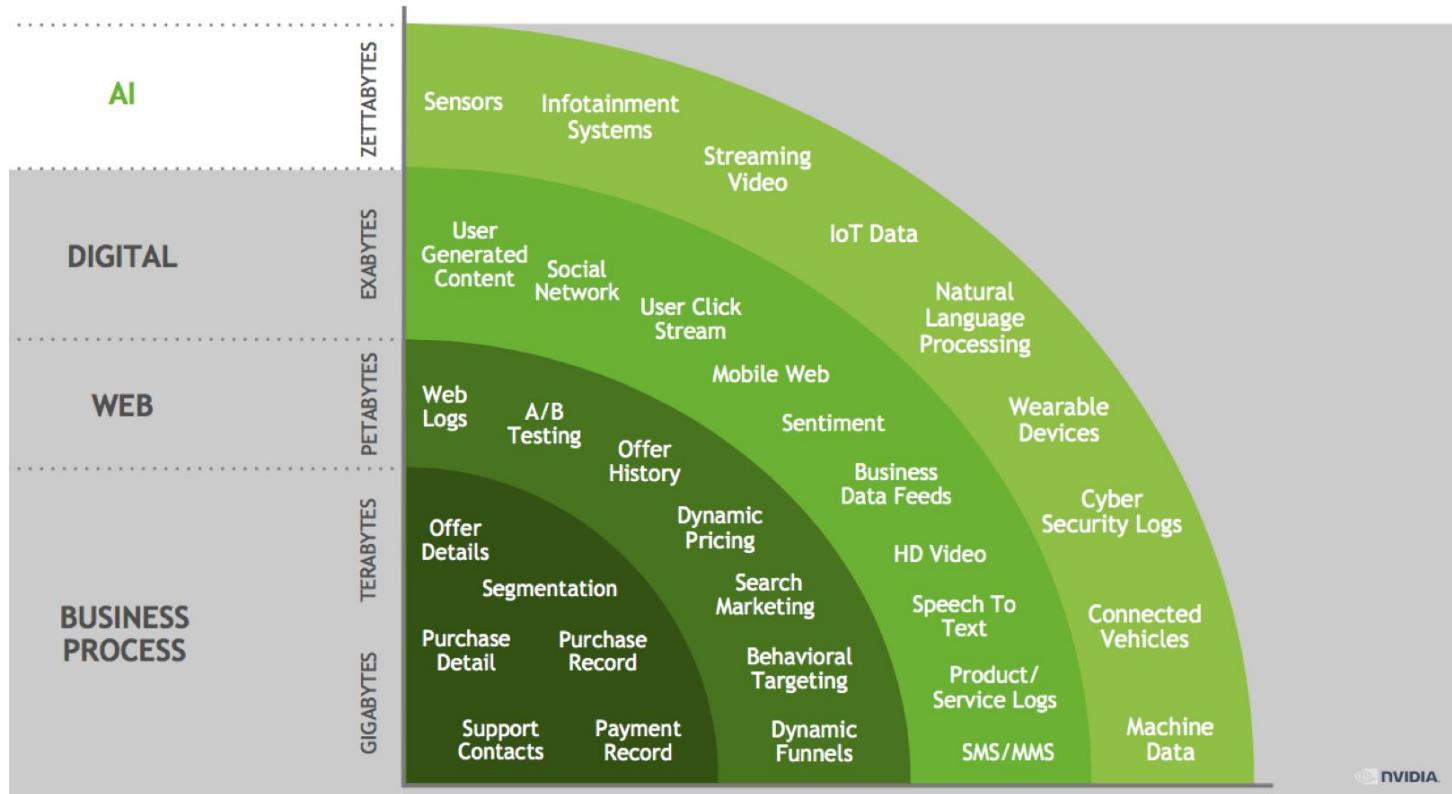


Active learning, semi-supervised learning, and weak supervision lead to dramatically more efficient solutions.

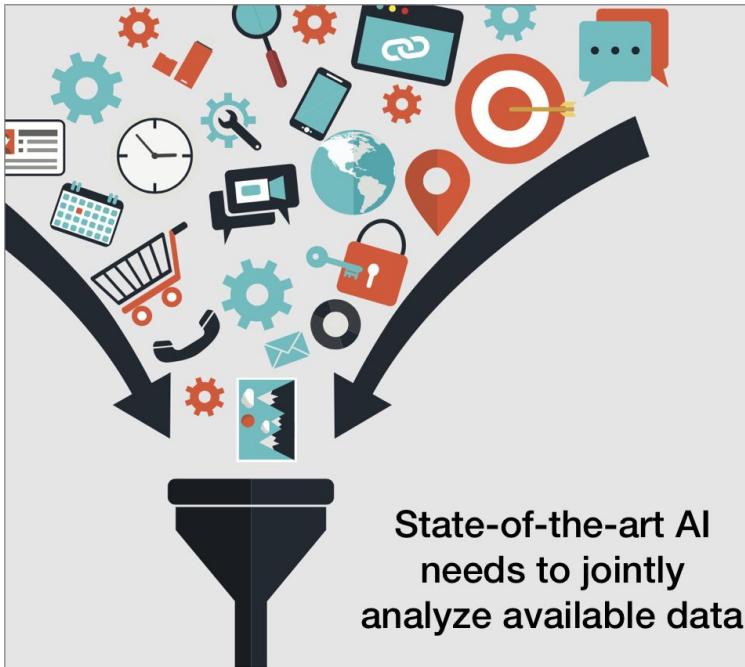
Outline

- 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



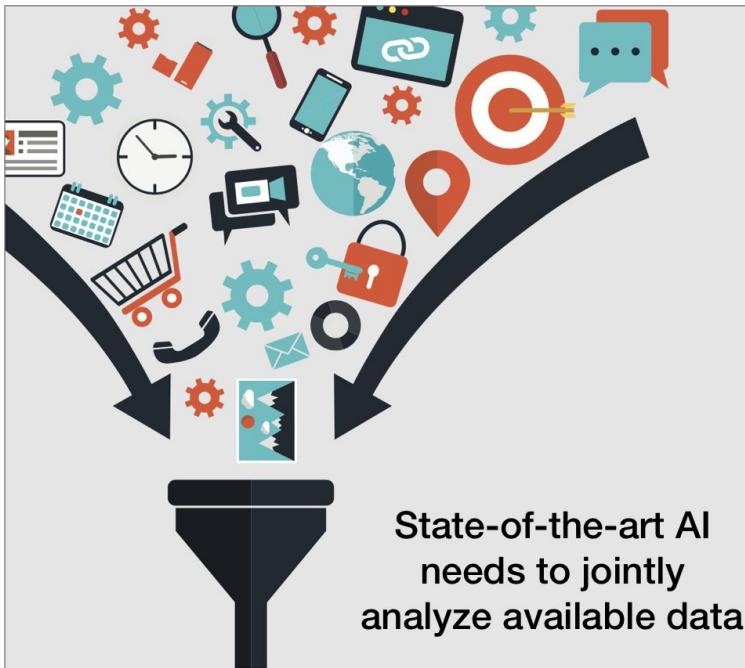
IMAGENET MovieLens



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



1990s (Features)

1970s (Rules)

Classical ML

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

Graphical models and logic

- Relational statistical learning
- Markov logic network



2009 (PGMs)

2010s

(Representation Learning)

Deep learning

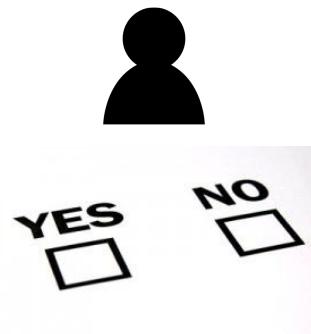
- Automatically learn representations
- Impressive with high-dimensional data
- Data hungry!

The ML Pipeline in the Deep Learning Era

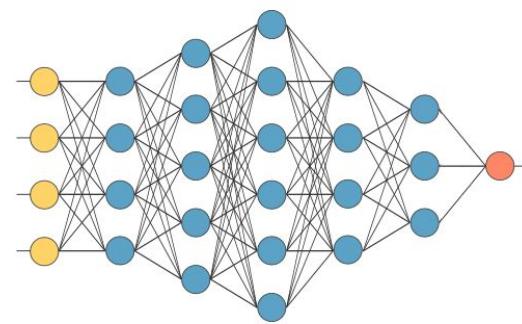
Data Collection



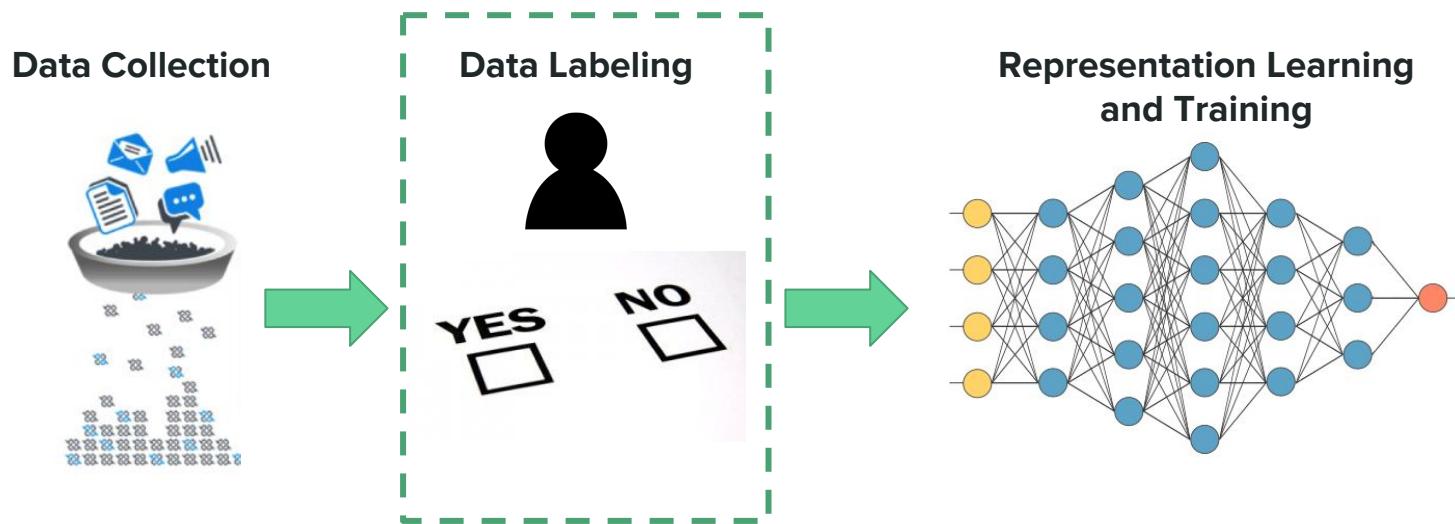
Data Labeling



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



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 20% 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)

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

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

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

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

Specificity (1 - false positive rate)

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

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

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

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

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

EM algorithm to obtain maximum-likelihood estimates.

Difference with Dawid-Skene is the estimation of 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).

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]

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]

Input: Labeling Functions,
Unlabeled data

DOMAIN
EXPERT

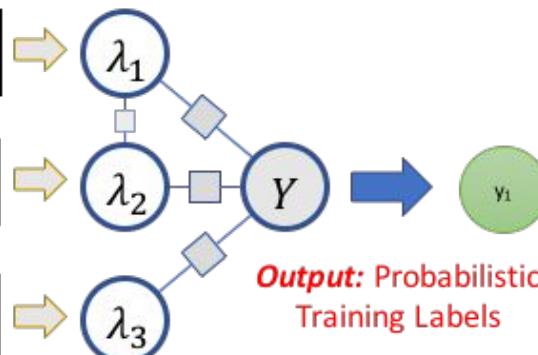


```
def lf1(x):
    cid = (x.chemical_id,
           x.disease_id)
    return 1 if cid in KB else 0
```

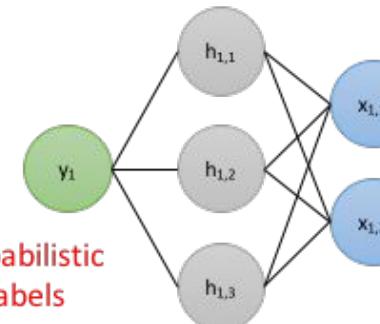
```
def lf2(x):
    m = re.search(r'.*cause.*',
                  x.between)
    return 1 if m else 0
```

```
def lf3(x):
    m = re.search(r'.*not
                  cause.*',
                  x.between)
    return 1 if m else 0
```

**Generative
Model**



**Noise-Aware
Discriminative Model**



*Ex. Application:
Knowledge Base
Creation (KBC)*



- 1 Users write *labeling functions* to generate noisy labels
- 2 We model the labeling functions' behavior to de-noise them
- 3 We use the resulting prob. labels to train a model

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

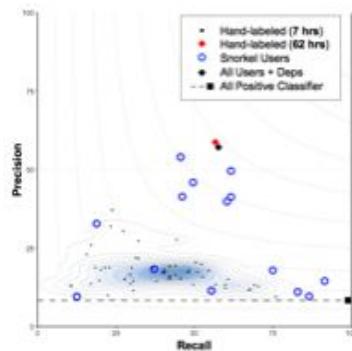


Snorkel biomedical workshop in collaboration with the NIH Mobilize Center



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



Marta Gaia Zanchi (@mzanchi) Following
For a newbie, I write pretty darn good #Snorkel #MachineLearning labeling functions. Thanks @MobilizeCenter @jasonafries @stevebach :)

3rd Place Score
No machine learning experience
Beginner-level Python

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

Alex Ratner



You can find Alex at the poster session tonight!

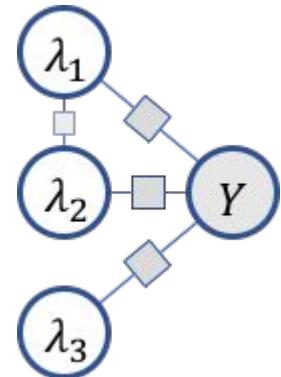
<https://ajratner.github.io>

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



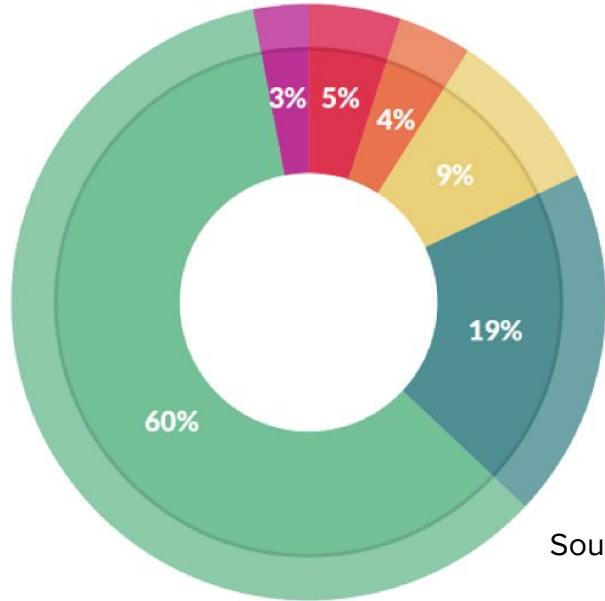
IMAGENET MovieLens



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%*

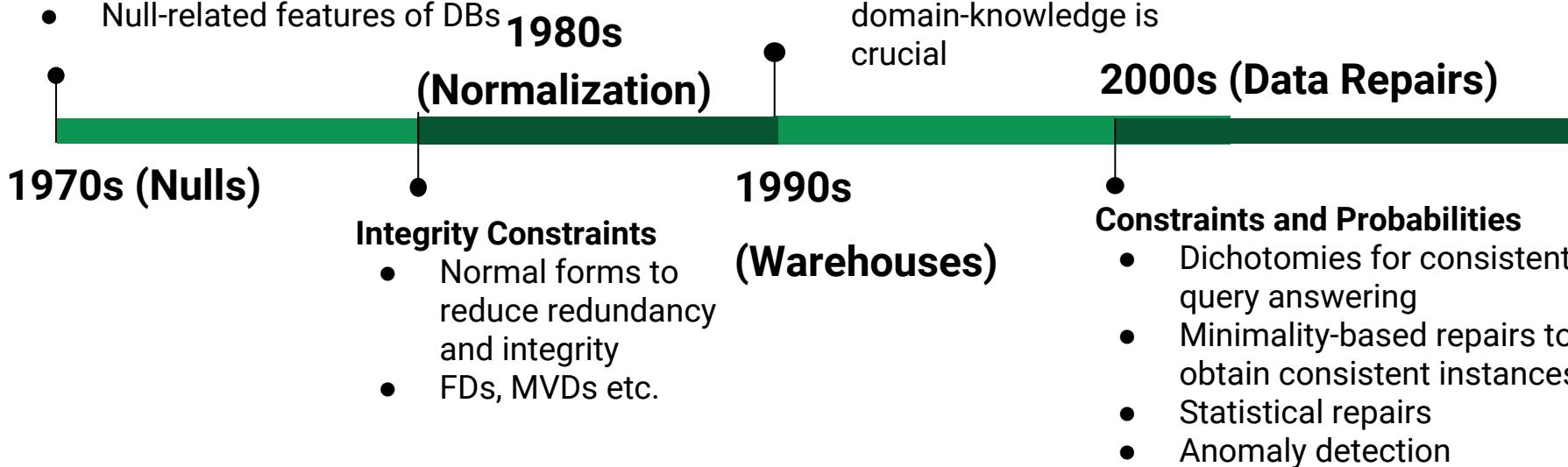
Source: Crowdflower

Cleaning and organizing data comprises 60% of the time spent on an analytics or 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

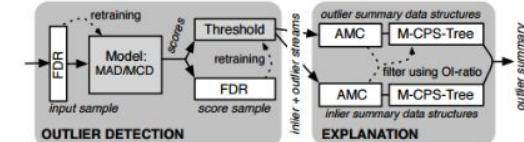
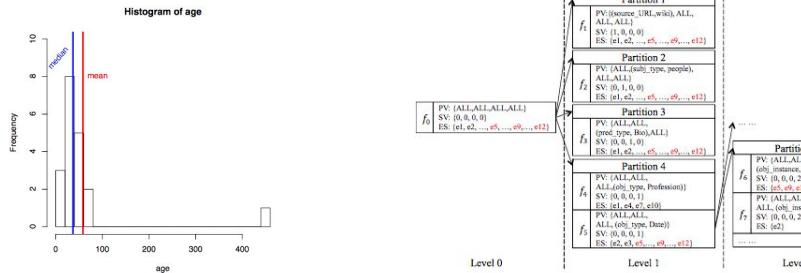


Where are we today?

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., SIGMOD'17]

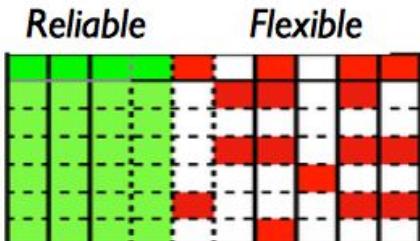


Where are we today?

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]



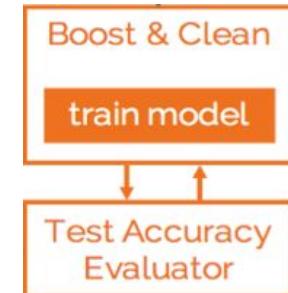
Each cell is a random variable

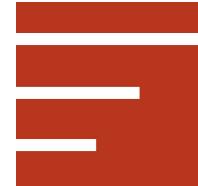
Constraints introduce correlations
c3: City, State, Address → Zip

External data introduce evidence

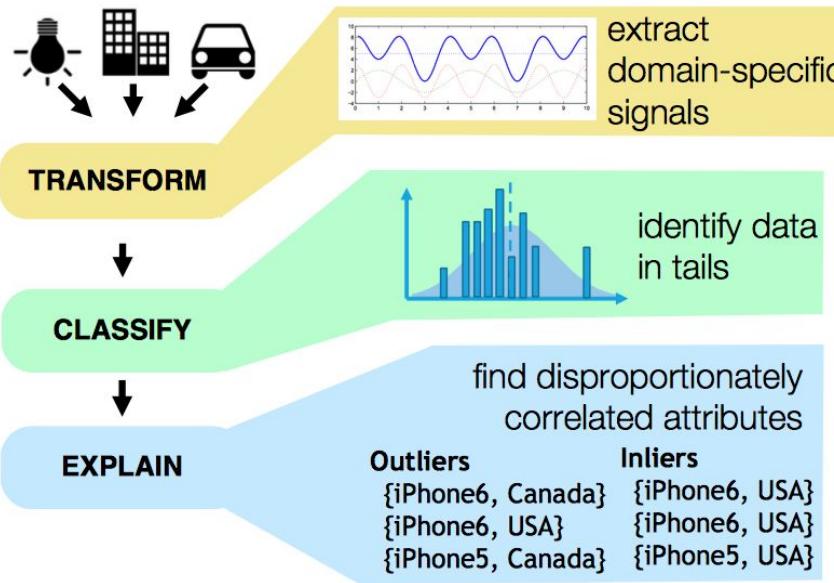
Address	City	State	Zip
3465 S Morgan ST	Chicago	IL	60608
3465 S Morgan ST	Chicago	IL	60609
3465 S Morgan ST	Chicago	IL	60609
3465 S Morgan ST	Chicago	IL	60608

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

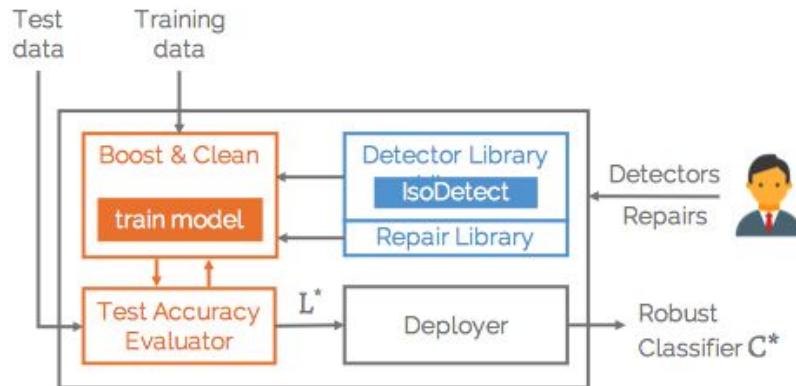
Sketch of a classifier for fast updates and queries for estimates of each weight and comes with approximation guarantees

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

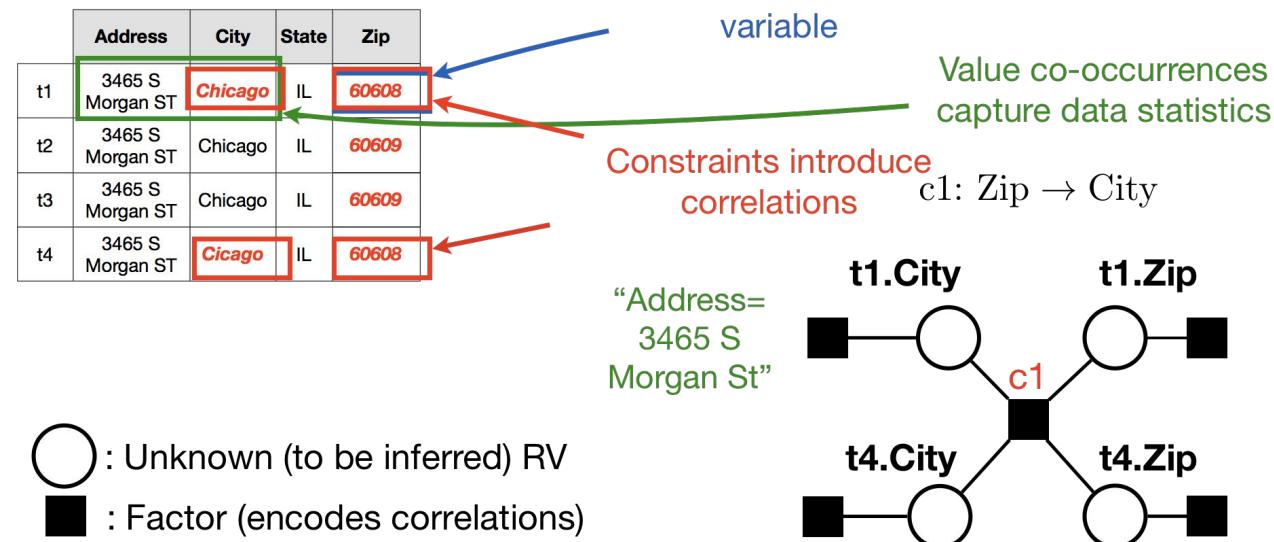


Code available at:

<http://www.holoclean.io>



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

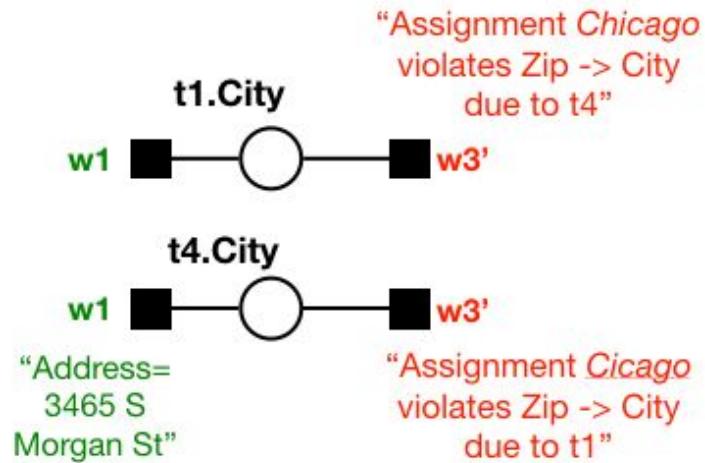


Holistic data cleaning framework: combines a variety of heterogeneous signals (e.g., integrity constraints, external knowledge, quantitative statistics)



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

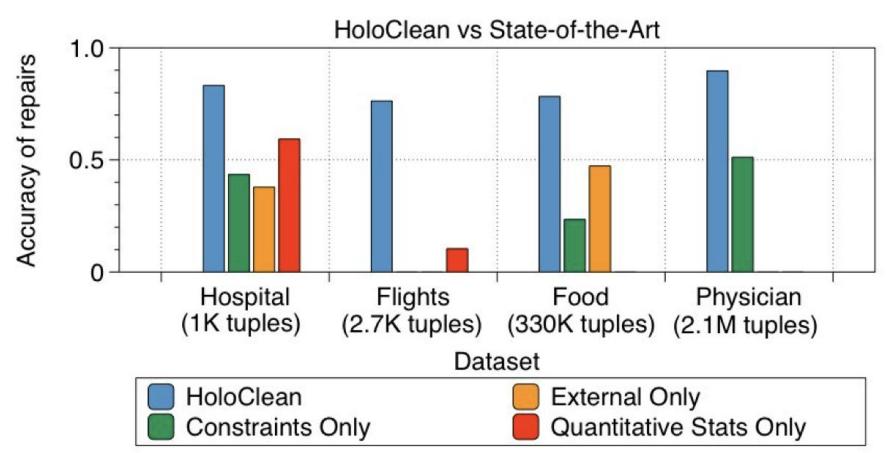
	Address	City	State	Zip
t1	3465 S Morgan ST	Chicago	IL	60608
t2	3465 S Morgan ST	Chicago	IL	60609
t3	3465 S Morgan ST	Chicago	IL	60609
t4	3465 S Morgan ST	Cicago	IL	60608



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

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

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

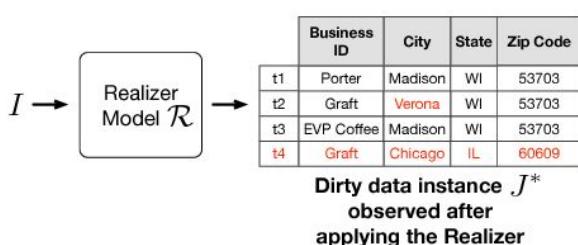
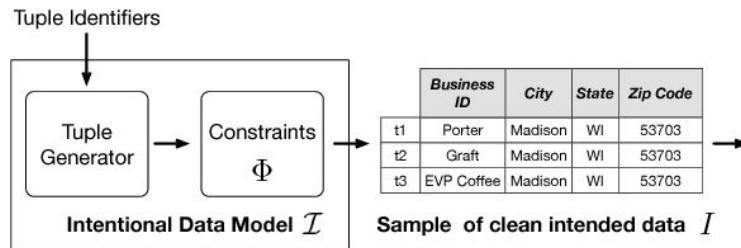
Probabilistic Unclean Databases [De Sa et al., 2018]

Unclean Database Generation

(A) Schema, Attribute Domain, and Constraint Specification

Tuple ID	Business Listing				Integrity Constraints	
Tuple Identifiers	Business ID	City	State	Zip Code	PK: Business ID FD: Zip Code → City, State	

(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

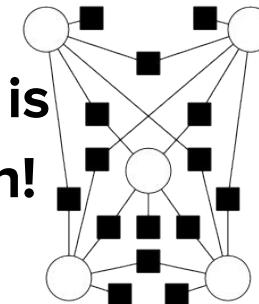
- Problem definition: **Detect and repair erroneous data.**
 - Short answers
 - **ML can help partly-automate cleaning**
Domain-expertise is still required.
 - **Scalability of ML-based data cleaning**
a pressing challenge. Exciting system!
 - **We need more end-to-end systems!**

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

Each cell is a random variable

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

External data introduce evidence



Outline

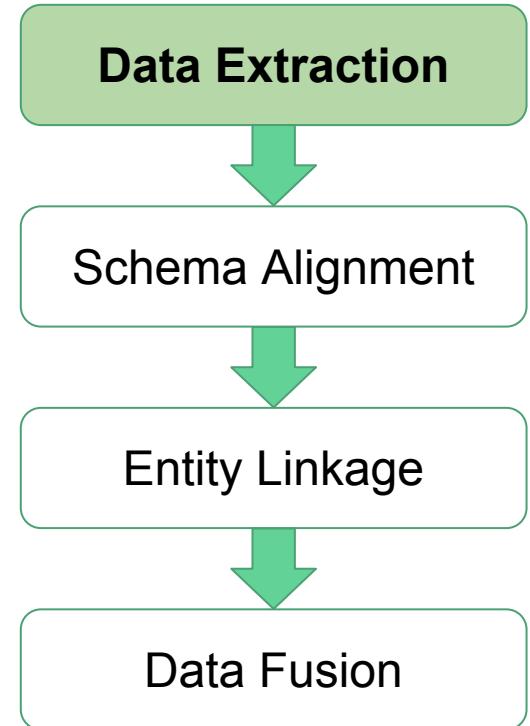
- Part I. Introduction
- Part II. ML for DI
- Part III. DI for ML
 - Creating training data
 - Data cleaning
- Part IV. Conclusions and research direction

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

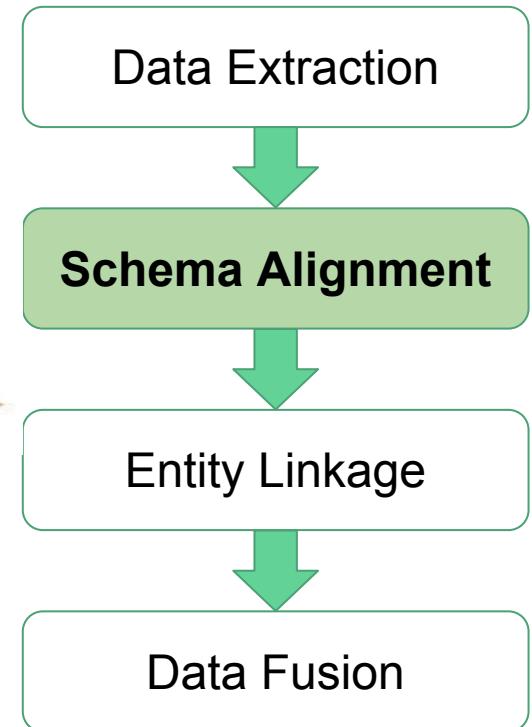
Revisit: Recipe for Data Extraction

- 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**
 - **DL effective for texts and LR is often effective for semi-stru data**



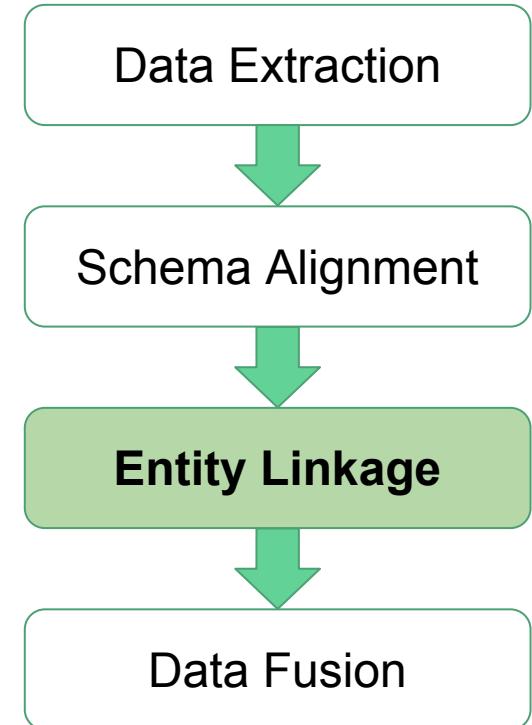
Revisit: Recipe for Schema Alignment

- Problem definition: **Align attributes with the same semantics**
- Short answers
 - **Interactive semi-automatic mapping**
 - **DL-based universal schema revived the field**
 - **Combine schema matching and universal schema for future**



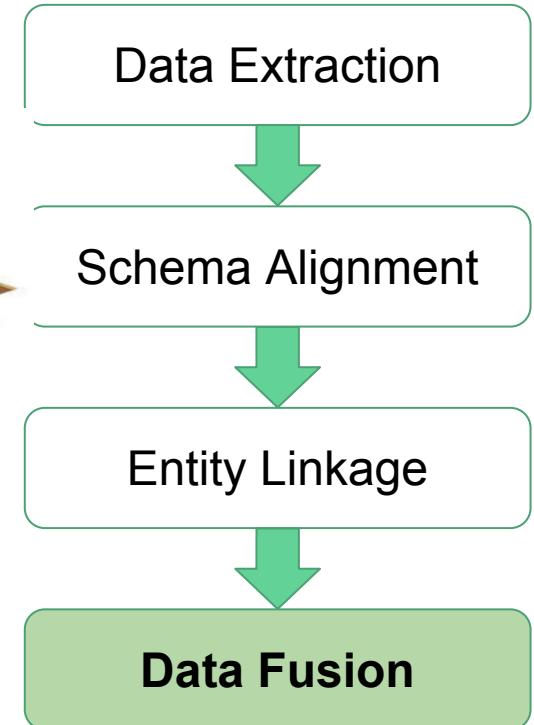
Revisit: Recipe for Entity Linkage

- 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**



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.

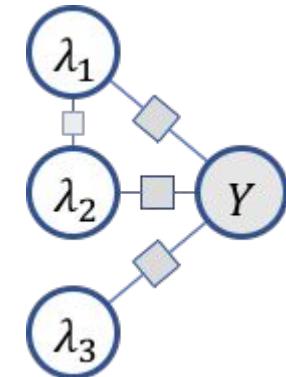


DI and ML: A Natural Synergy

- 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

Revisit: 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.



Recipe for Data Cleaning

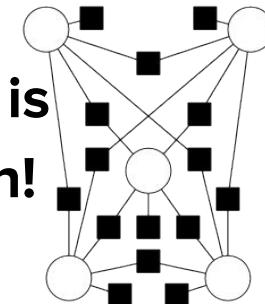
- Problem definition: Detect and repair erroneous data.
 - Short answers
 - ML can help partly-automate cleaning
Domain-expertise is still required.
 - Scalability of ML-based data cleaning a pressing challenge. Exciting systems!
 - We need more end-to-end systems!

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

Each cell is a random variable

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

External data introduce evidence



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 we 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



MLPerf

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. *VLDB*, 11(3), pp.269–282.
- Rekatsinas, T. et al., 2017. HoloClean: Holistic Data Repairs with Probabilistic Inference. *VLDB*, 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*.
- Veroios, 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. Fonduer: 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. Newslett.*, 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.