

Entity Linkage with Deep Learning

Shobeir Fakhraei
University of Southern California

Where are we?

- Information Extraction
 - Semi-structured: Web, HTML
 - Unstructured: Text, Ads, Tweets, ...
- Entity Linkage, Data Cleaning, Normalization
- Logical Data Integration
 - Mediators, Query Rewriting
 - Warehouse, Logical Data Exchange
- Automatic Source Modeling/Learning Schema Mappings
- Semantic Web
 - RDF, SPARQL, OWL, Linked Data
- Advanced Topics
 - Geospatial Data Integration, Knowledge Graphs

String Matching

...

Blocking

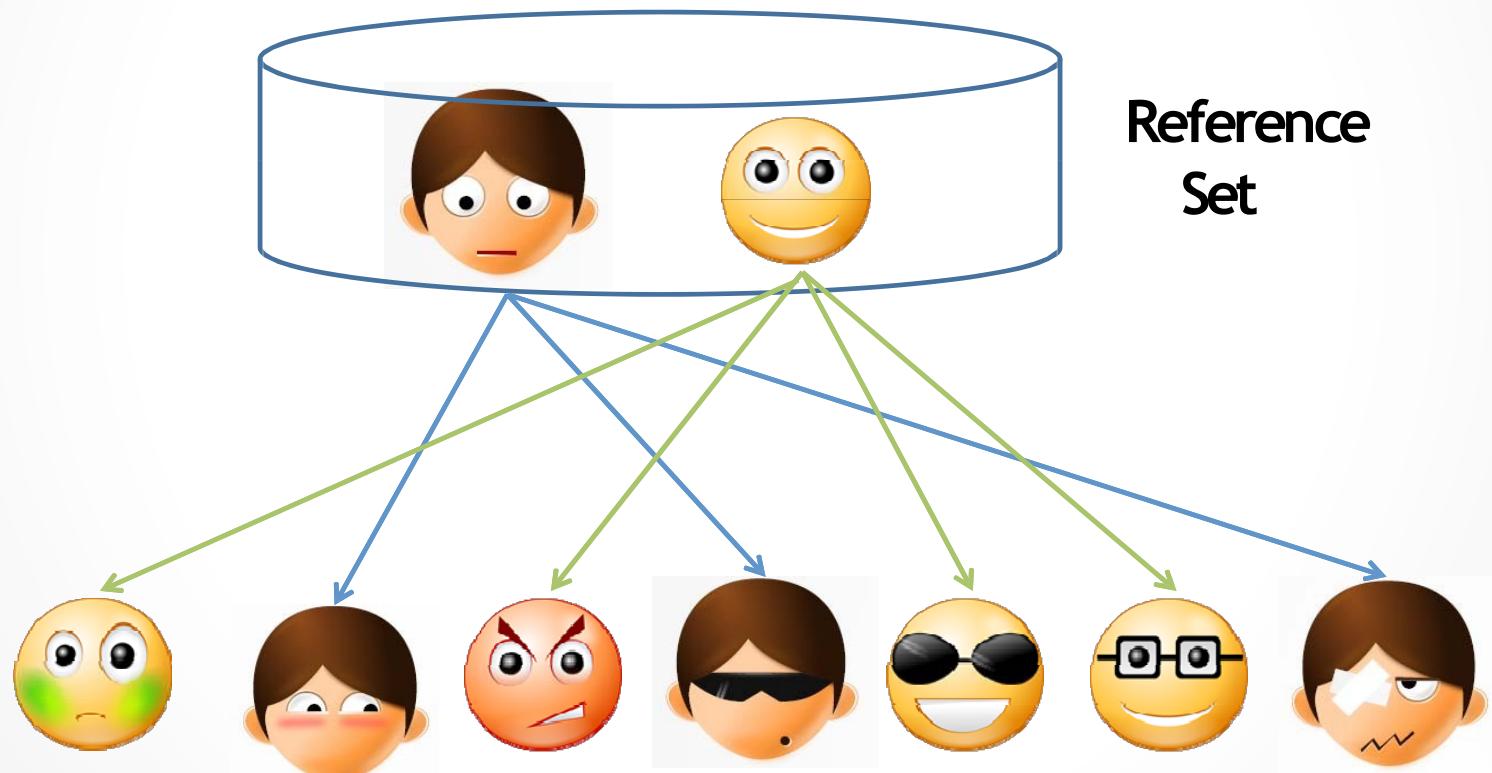


Deep Entity Normalization

...

Reference Matching Problem

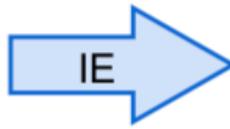
- Match noisy records to clean records in a reference table



Biomedical Information Extraction



textual abstract:
summary for human



Subject	Relation	Object
p53	is_a	protein
Bax	is_a	protein
p53	has_function	apoptosis
Bax	has_function	induction
apoptosis	involved_in	cell_death
Bax	is_in	mitochondrial outer membrane
Bax	is_in	cytoplasm
apoptosis	related_to	caspase activation
...

structured knowledge extraction: summary for machine

Biomedical Information Extraction

The figure illustrates the process of Biomedical Information Extraction (BIE) from a scientific publication. It shows a flow from the original document to a textual abstract for humans and then to structured data representation.

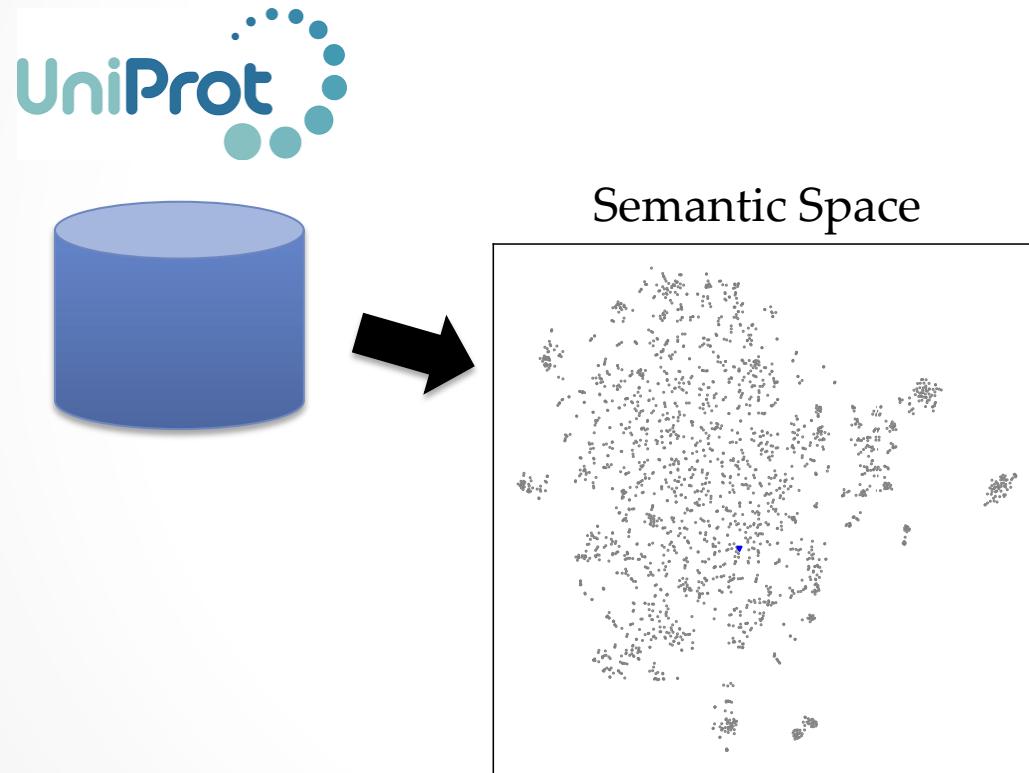
Textual Abstract: A red box highlights a portion of the abstract from a paper titled "Involvement of Tumor Necrosis Factor Receptor-associated Protein 1 (TRAP1) in Apoptosis Induced by β -Hydroxyisovalerylshikimic". The abstract discusses the role of TRAP1 in apoptosis induced by β -Hydroxyisovalerylshikimic. A red arrow points from this section to the text "textual abstract: summary for human".

UniProtKB Search Results: A screenshot of the UniProtKB interface for the protein P04637 (P53_HUMAN). The search bar shows "p53". The results page displays the protein's name, TP53, and its function as a cellular tumor antigen. A red arrow points from the UniProtKB search results to a table showing extracted triples.

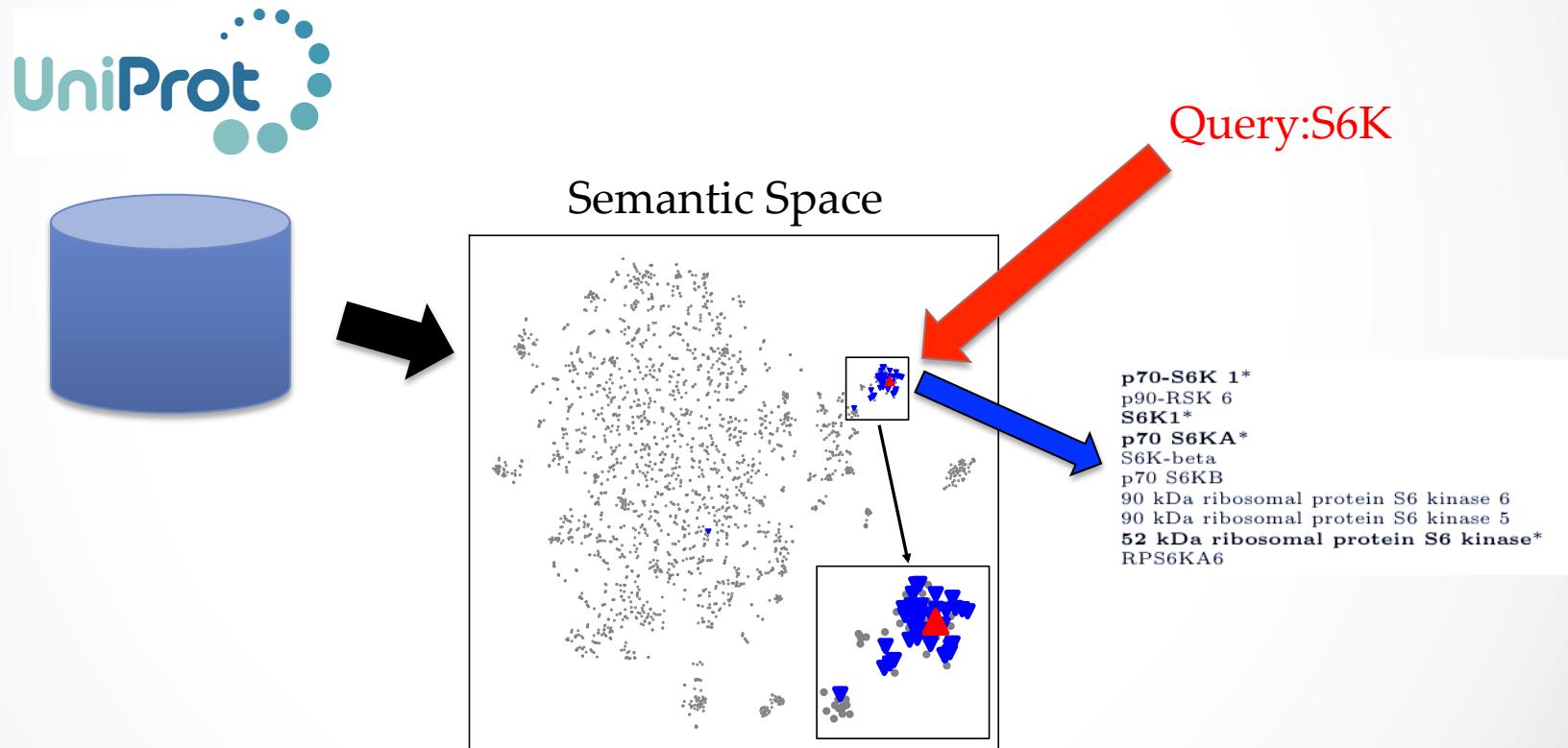
Extracted Triples: A table showing the extracted triples from the UniProtKB search results. The triples are:

Subject	Relation	Object
p53	is_a	protein
Bax	is_a	protein

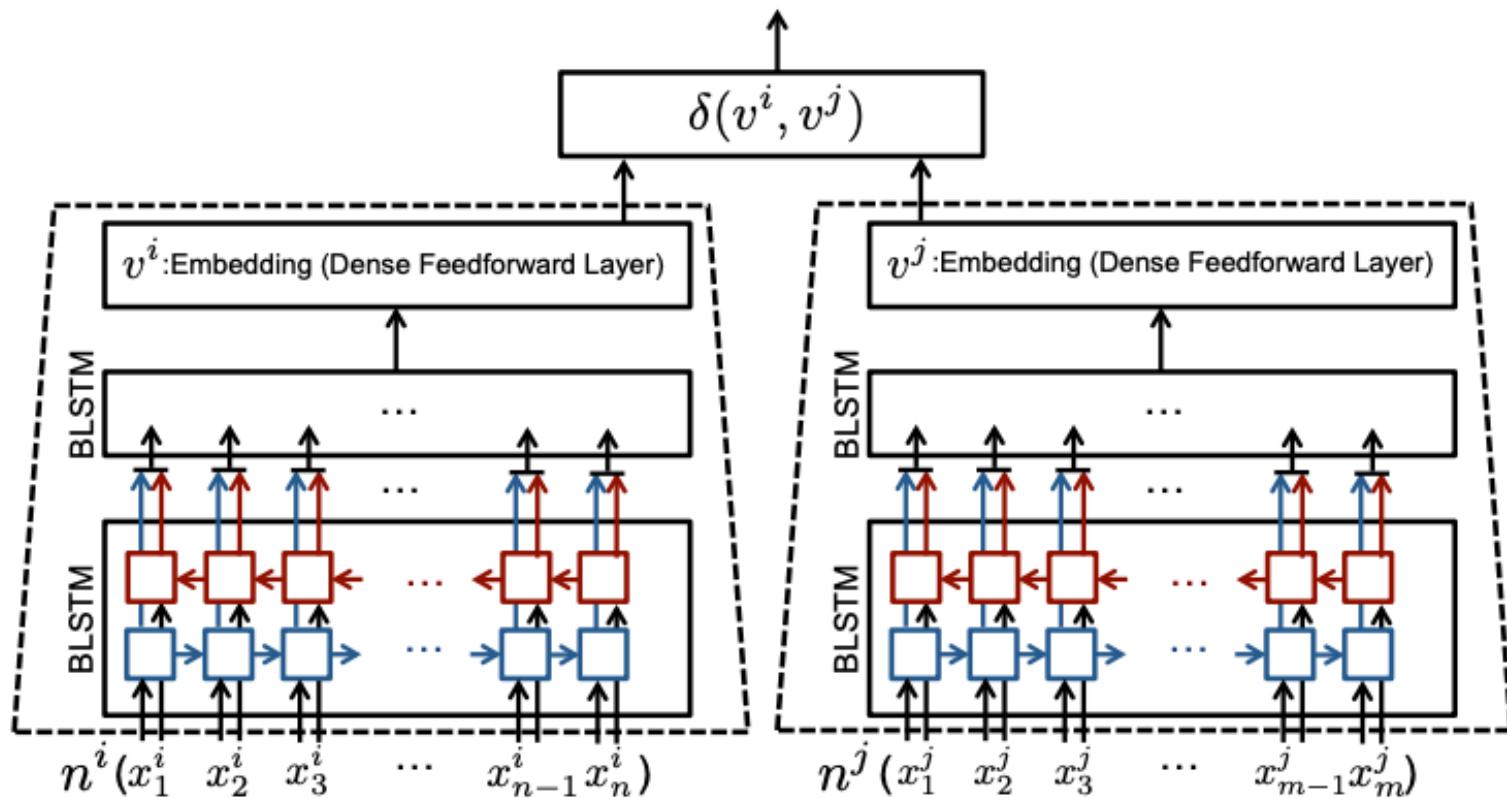
Our Deep Learning-based Solution



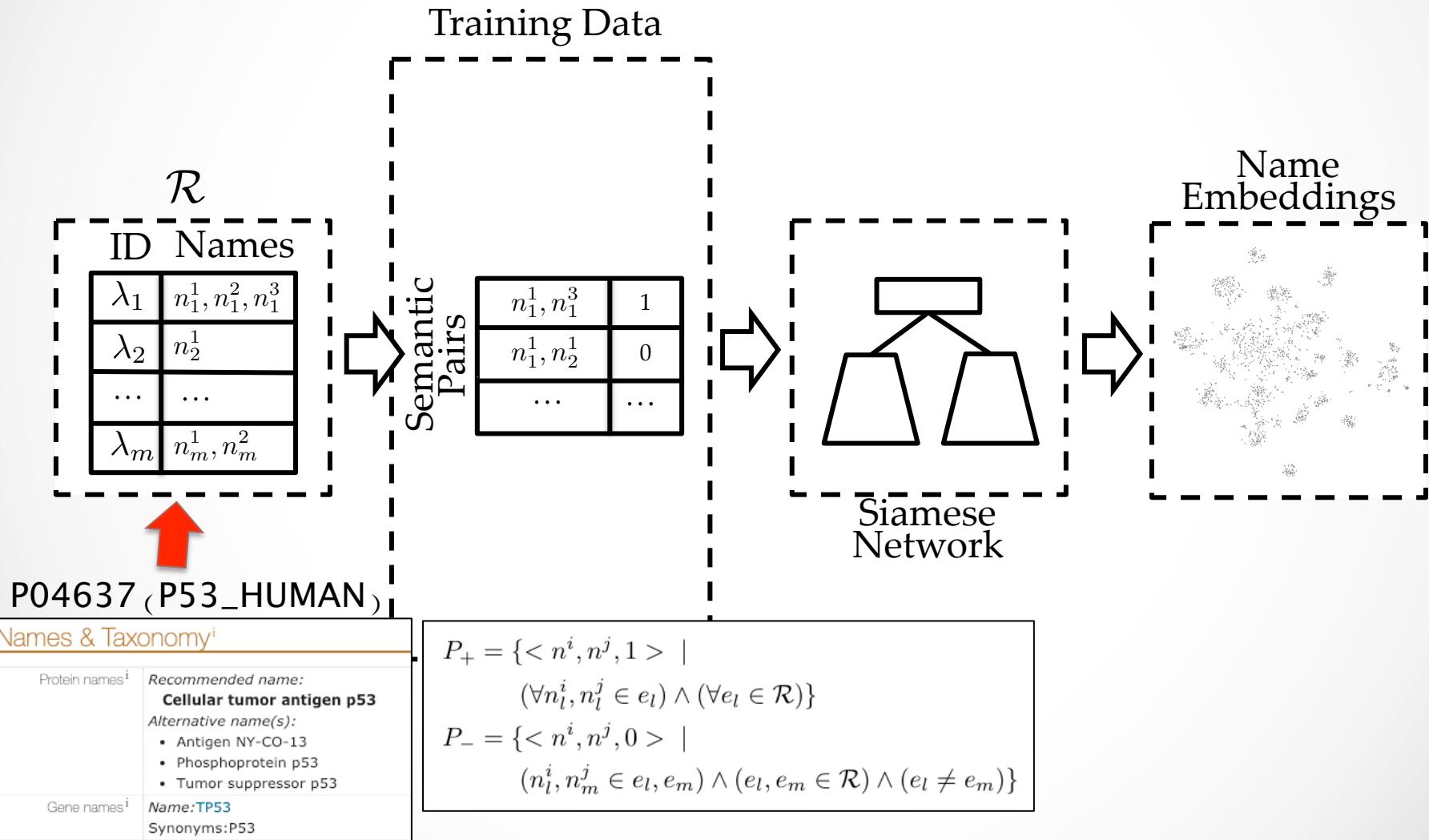
Our Deep Learning-based Solution



Similarity Learning via Siamese Network

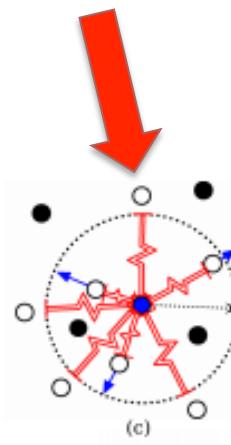
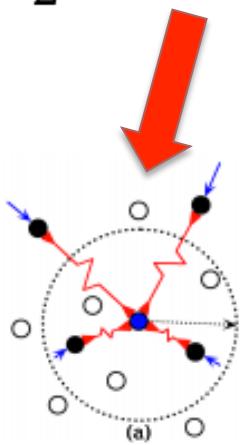


Pair Generation: Semantic

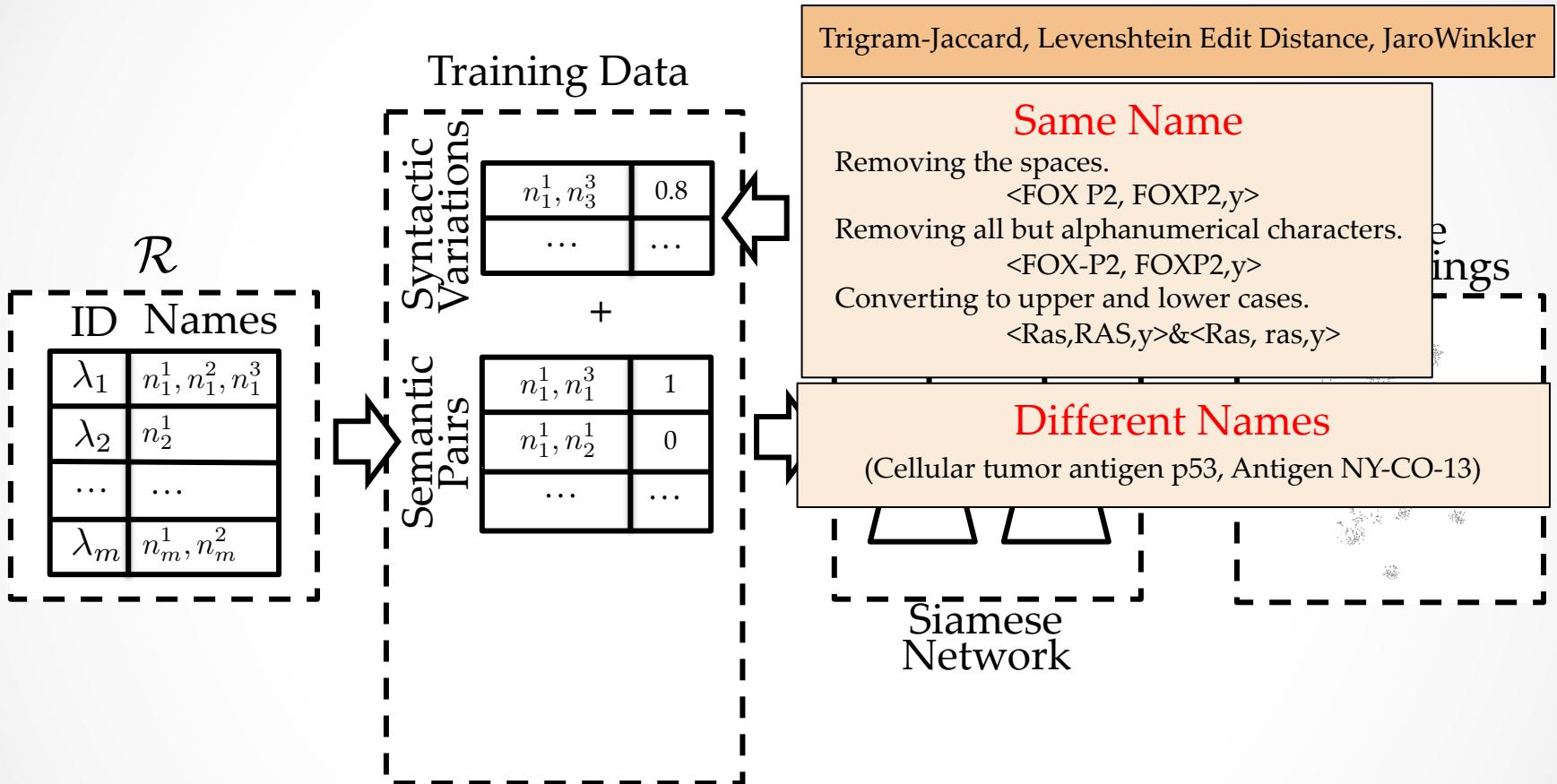


Contrastive Loss

$$\mathcal{L} = \frac{1}{2}y\delta(v_i, v_j)^2 + \frac{1}{2}(1-y)\max(0, m - \delta(v_i, v_j))^2$$

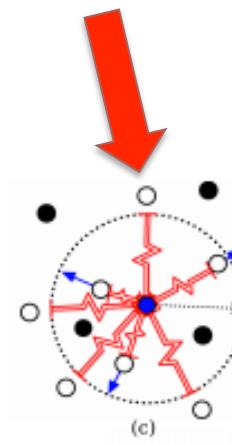
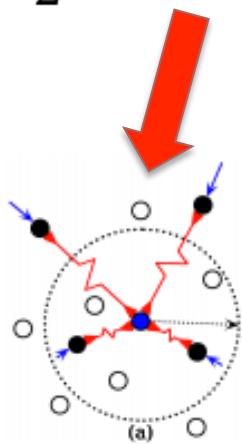


Pair Generation: Syntactic Variations



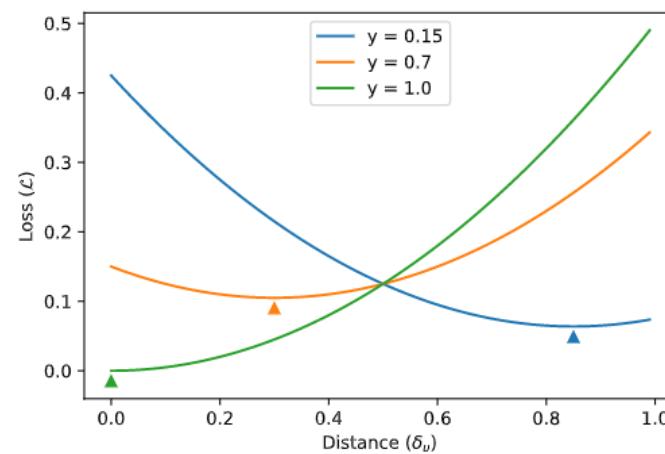
Contrastive Loss with soft labels

$$\mathcal{L} = \frac{1}{2}y\delta(v_i, v_j)^2 + \frac{1}{2}(1-y)\max(0, m - \delta(v_i, v_j))^2$$

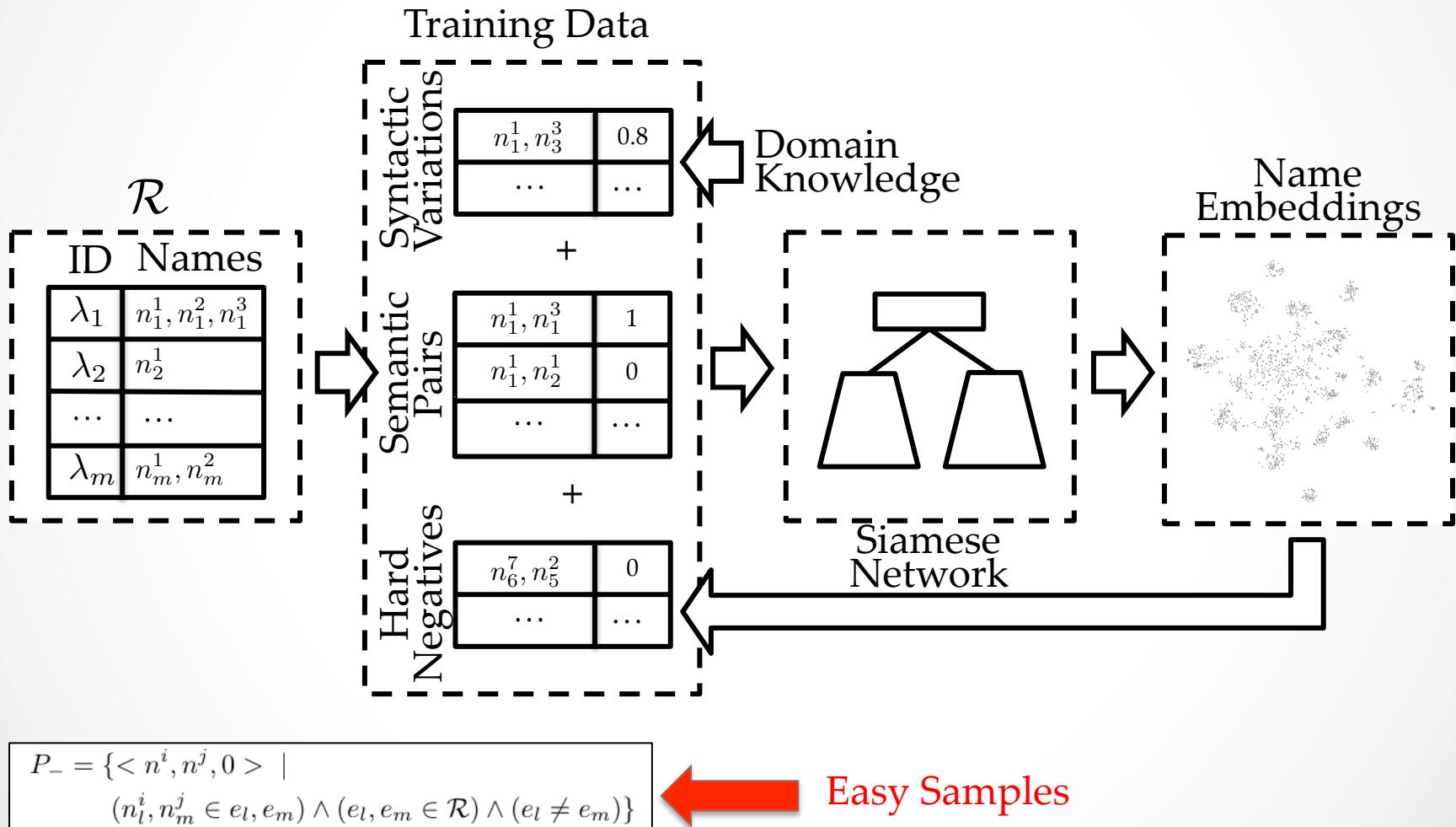


$$\frac{\partial \mathcal{L}}{\partial \delta_v} = y\delta_v - (1-y)(1-\delta_v)$$

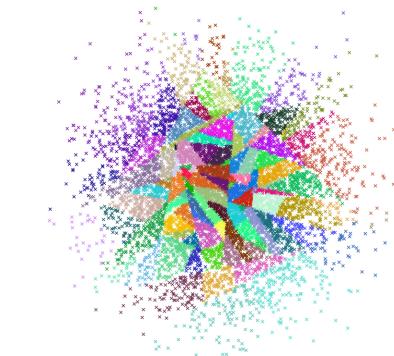
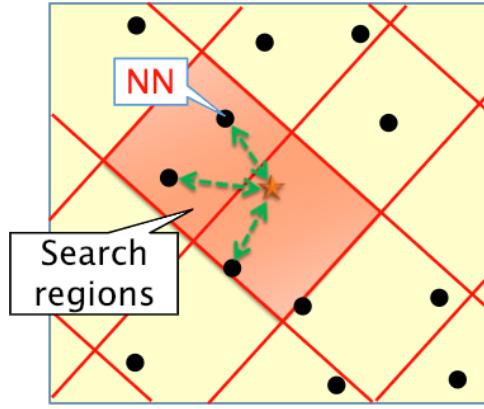
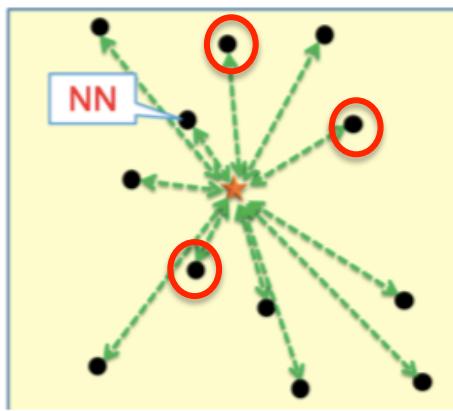
$$\arg \min_{\delta_v} \mathcal{L} = \{\delta_v \mid y + \delta_v - 1 = 0\} = 1 - y$$



Pair Generation: Hard Negatives

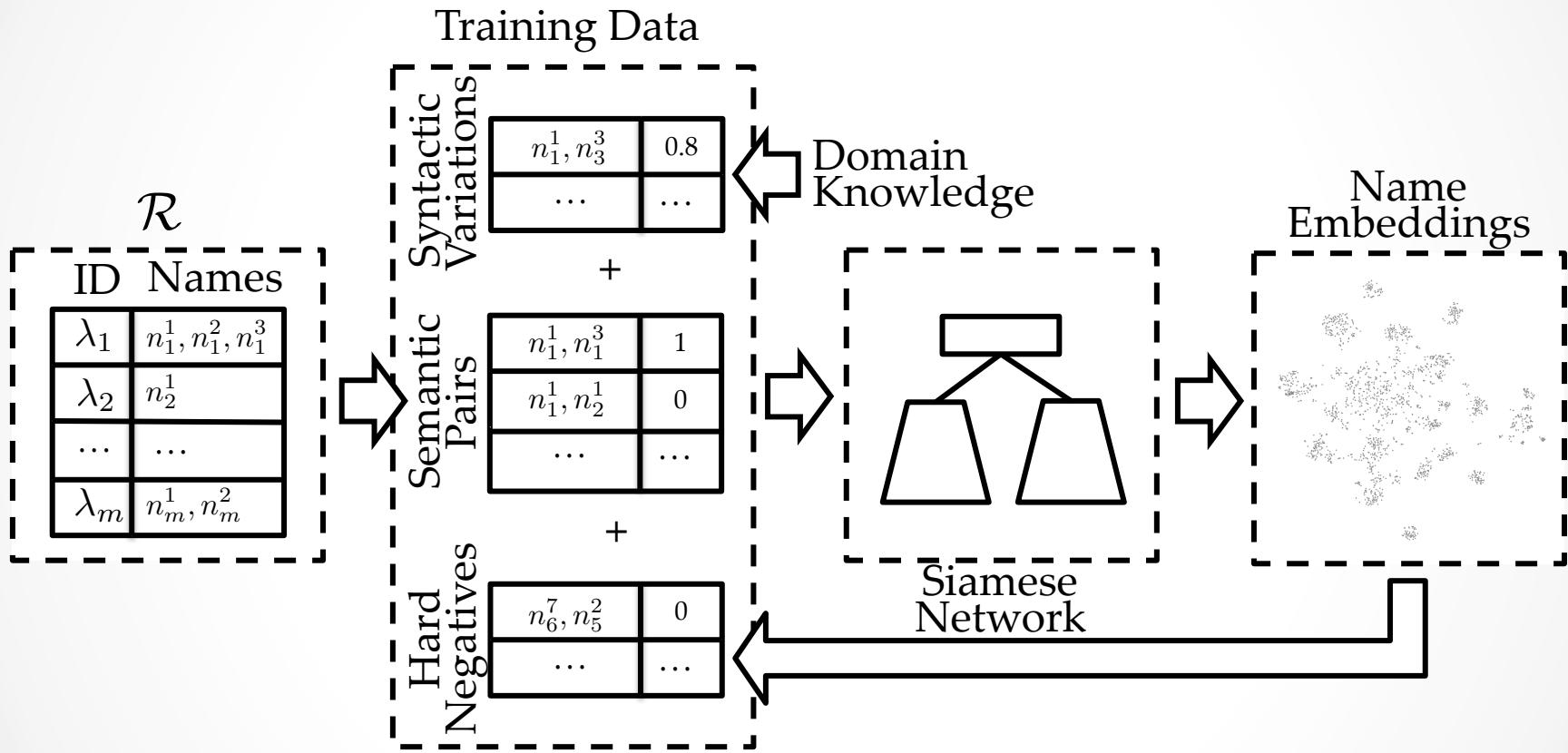


Finding Hard Negatives



Annoy
(Approximate Nearest
Neighbors Oh Yeah)

Overall Framework (Similarity Learning)



Similarity Learning

Algorithm 1 NSEEN: Similarity Learning

- 1: **procedure** TRAINSIM(\mathcal{R}, P_d)
- 2: **Input:** \mathcal{R} reference set
- 3: **Input:** P_d pairs based on knowledge of syntactic variation in the domain
- 4: Generate pairs based on reference set \mathcal{R} and add them to training data \mathcal{D}
- 5: Add P_d pairs to the training data \mathcal{D}
- 6: **for** k times **do**
- 7: Train the model \mathcal{M} (Siamese network) on \mathcal{D}
- 8: Embed all the names in \mathcal{R} : $n \rightarrow v$
- 9: **for all** v_l^i **do** ▷ Hard negative mining
- 10: find the k closest v_k^j to v_l^i
- 11: **if** $k \neq l$ **then**
- 12: add $< n_k^j, n_l^i, 0 >$ to training data \mathcal{D}
- 13: **return** \mathcal{M} ▷ The trained embedding function

Embedding the Database

Algorithm 2 NSEEN: Embedding & Hashing \mathcal{R}

```
1: procedure EMBED( $\mathcal{R}, \mathcal{M}$ )
2:   for all  $n_i \in \mathcal{R}$  do
3:      $n_i \xrightarrow{\mathcal{M}} v_i$ 
4:   for all  $v_i$  do
5:     Hash  $v_i$  and store in  $\mathcal{H}_{v_i}$ 
6:   return  $\mathcal{H}_v$                                  $\triangleright$  Hashed embeddings
```

Retrieval

Algorithm 3 NSEEN: Retrieval

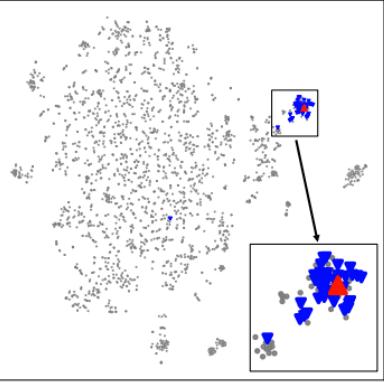
- 1: **procedure** RETRIEVE(\mathcal{H}_v , \mathcal{M} , n_q)
 - 2: Embed the query name: $n_q \xrightarrow{\mathcal{M}} v_q$
 - 3: Find the closest v_k^j to v_q using approximate nearest neighbor search (Annoy) on \mathcal{H}_v
 - 4: **return** λ_k as the ID (i.e., λ_q)
-

Results

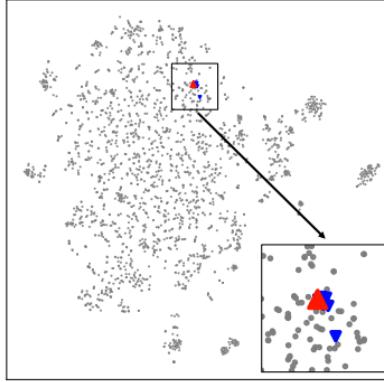
\mathcal{R}	DS	Model	H@1	H@3	H@5	H@10
UniProt	BC1	NSEEN	0.833	0.869	0.886	0.894
		Baseline	0.814	0.864	0.875	0.885
	BC2	NSEEN	0.861	0.888	0.904	0.930
		Baseline	0.841	0.888	0.904	0.919
ChEBI	BC1	NSEEN	0.505	0.537	0.554	0.574
		Baseline	0.418	0.451	0.460	0.468
	BC2	NSEEN	0.578	0.608	0.624	0.641
		Baseline	0.444	0.472	0.480	0.491

Examples

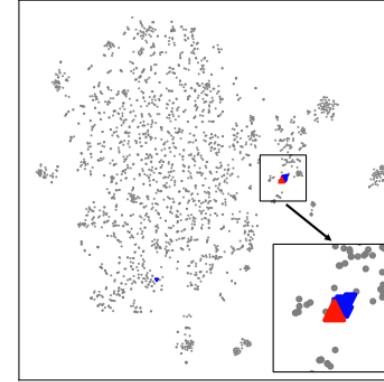
S6K	PLC γ 2	IKK ϵ	H3
<p>p70-S6K 1*</p> <p>p90-RSK 6</p> <p>S6K1*</p> <p>p70 S6KA*</p> <p>S6K-beta</p> <p>p70 S6KB</p> <p>90 kDa ribosomal protein S6 kinase 6</p> <p>90 kDa ribosomal protein S6 kinase 5</p> <p>52 kDa ribosomal protein S6 kinase*</p> <p>RPS6KA6</p>	<p>PLC-gamma-2*</p> <p>PLC-gamma-1</p> <p>PLCG2*</p> <p>Phospholipase C-gamma-2*</p> <p>Phospholipase C-gamma-1</p> <p>PLC</p> <p>PLCG1</p> <p>Phosphoinositide phospholipase C-gamma-2*</p> <p>PLC-IV*</p> <p>PLCB</p>	<p>IKK-epsilon*</p> <p>IKKE*</p> <p>I-kappa-B kinase epsilon*</p> <p>IkBKE*</p> <p>IKBKE*</p> <p>IKBE</p> <p>IK1</p> <p>IK1</p> <p>IKKG</p> <p>INKA1</p>	<p>Histone H3/a*</p> <p>Histone H3/o*</p> <p>Histone H3/m*</p> <p>Histone H3/b*</p> <p>Histone H3/f*</p> <p>HIST1H3C*</p> <p>Histone H3/k*</p> <p>Histone H3/i*</p> <p>HIST1H3G*</p> <p>Histone H3/d*</p>



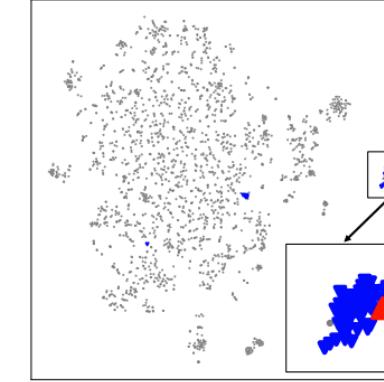
(a) S6K



(b) PLC γ 2



(c) IKK ϵ



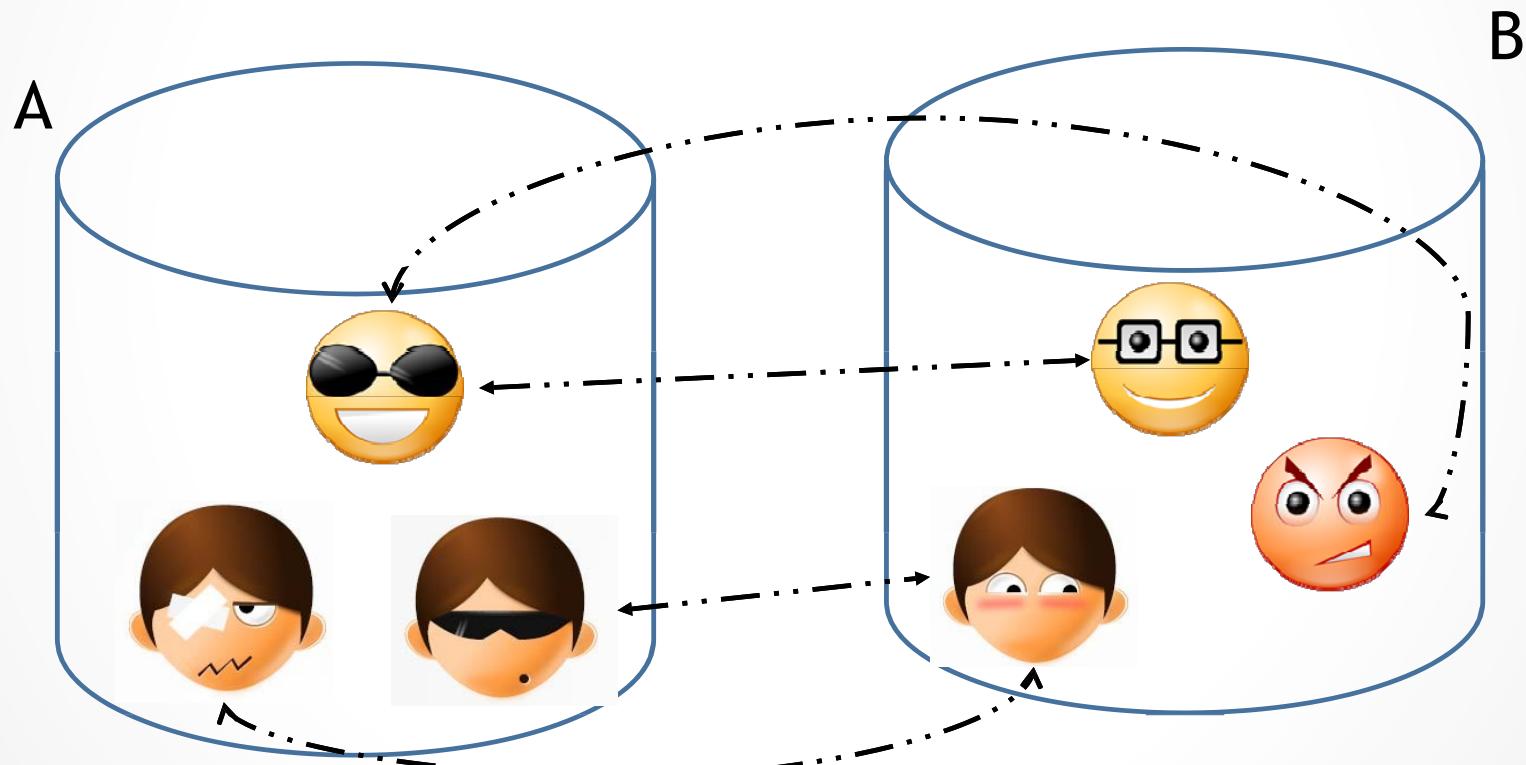
(d) H3

Deep Record Matching

...

Record Linkage Problem Statement

- Link records that match across databases



Different Kinds of Records

	Name	City	Age
t_1	Dave Smith	New York	18

	Name	City	Age
t_2	David Smith	New York	18

(a) structured

	Name	Brand	Price
t_1	Adobe Acrobat 8		299.99

	Name	Brand	Price
t_2	Acrobat 8	Adobe	299.99

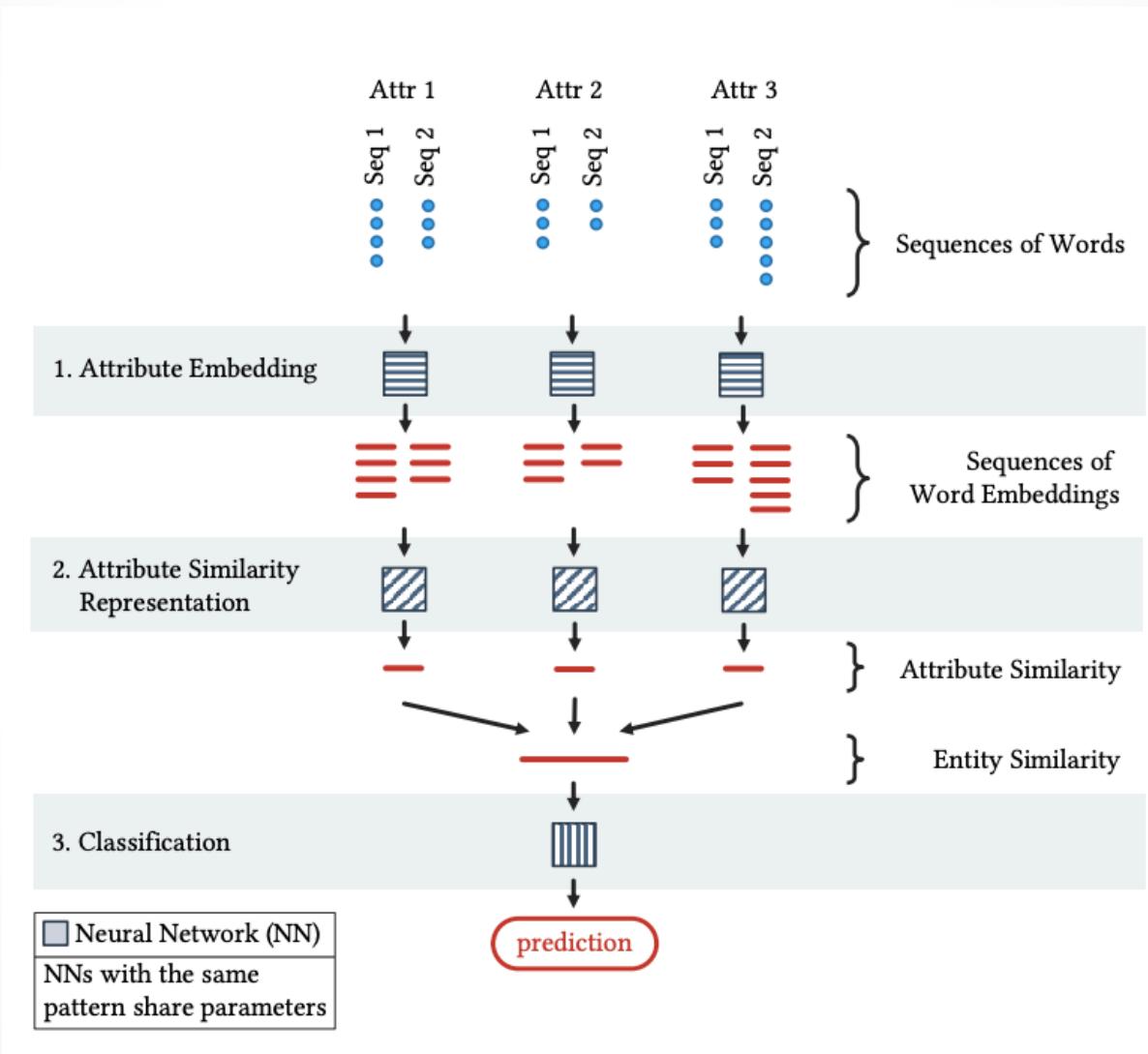
(c) dirty

	Description
t_1	Kingston 133x high-speed 4GB compact flash card ts4gcf133, 21.5 MB per sec data transfer rate, dual-channel support, multi-platform compatibility.

	Description
t_2	Kingston ts4gcf133 4GB compactflash memory card (133x).

(b) textual

Overall Architecture



Options at each stage

Architecture module	Options	
Attribute embedding	<i>Granularity:</i> (1) Word-based (2) Character-based	<i>Training:</i> (3) Pre-trained (4) Learned

Figure 3: The design space of DL solutions for EM.

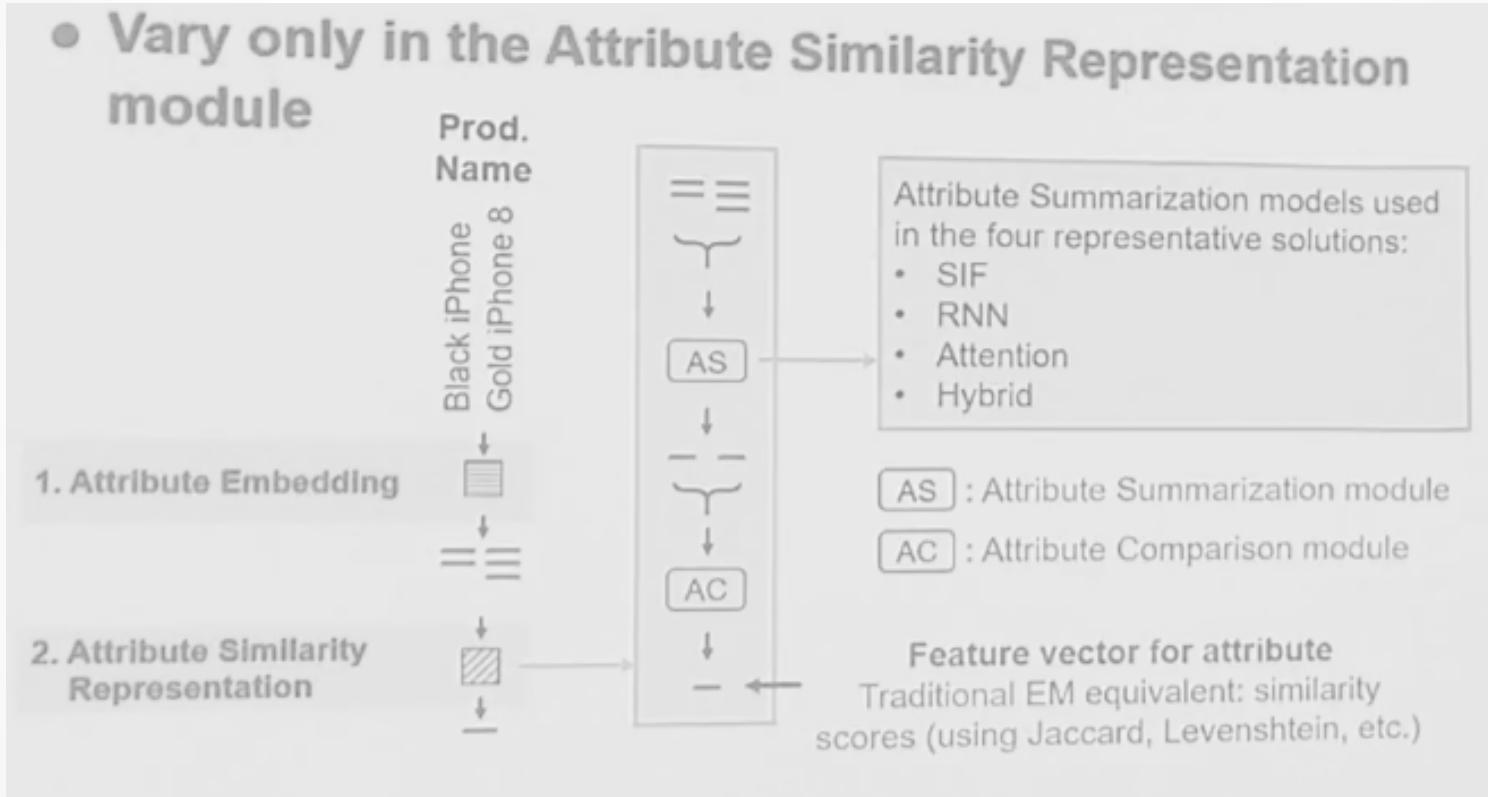
Options at each stage

Architecture module		Options	
Attribute embedding		<i>Granularity:</i> (1) Word-based (2) Character-based	<i>Training:</i> (3) Pre-trained (4) Learned
Attribute similarity representation	(1) Attribute summarization	(1) Heuristic-based (2) RNN-based (3) Attention-based (4) Hybrid	
	(2) Attribute comparison	(1) Fixed distance (cosine, Euclidean) (2) Learnable distance (concatenation, element-wise absolute difference, element-wise multiplication)	

Figure 3: The design space of DL solutions for EM.

Attribute Similarity

- Vary only in the Attribute Similarity Representation module

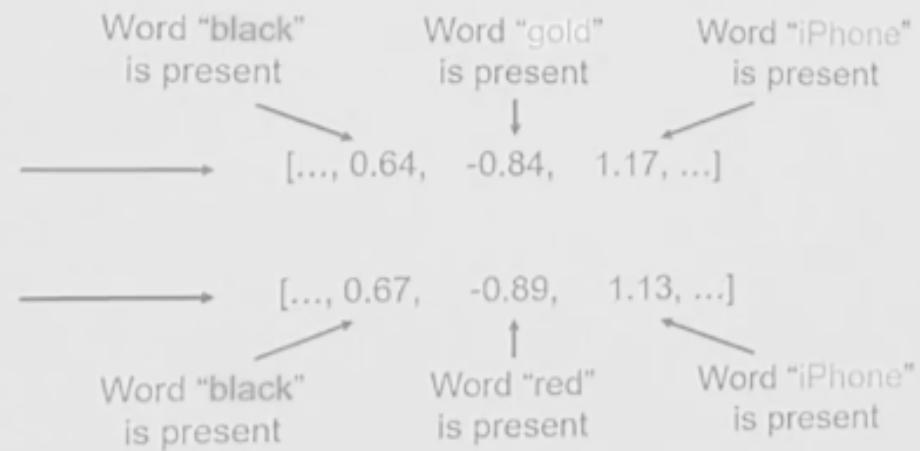


Aggregation-based attribute summarization

- Vector captures which words are present
 - A “bag of words” representation

Product Name
black iPhone with gold case

Product Name
black iPhone with a red case



SIF: Attribute Aggregation

- Uses smooth inverse frequency
- Weighted sum of word embeddings
 - Weight is inversely proportional to freq. of word

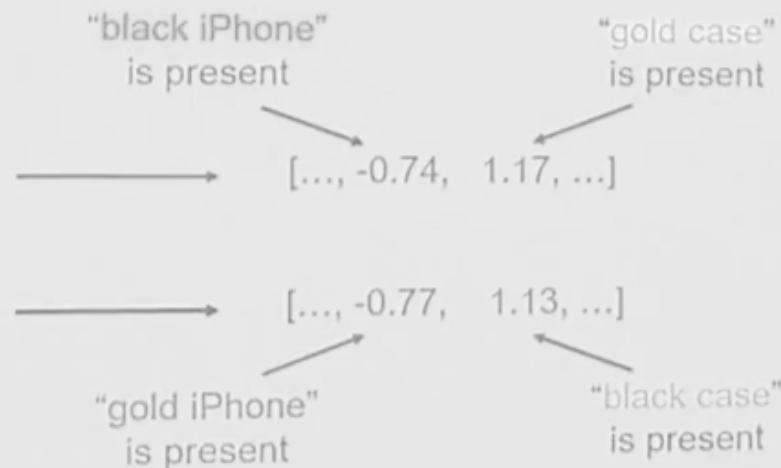
	Tuple pair	Word weights	SIF attribute summary computation													
t_1	<table border="1"><thead><tr><th>Product Name</th></tr></thead><tbody><tr><td>Laptop HP191</td></tr></tbody></table>	Product Name	Laptop HP191	<table border="1"><thead><tr><th>Word</th><th>Weight</th></tr></thead><tbody><tr><td>Laptop</td><td>0.05</td></tr><tr><td>Notebook</td><td>0.07</td></tr><tr><td>HP191</td><td>0.9</td></tr></tbody></table>	Word	Weight	Laptop	0.05	Notebook	0.07	HP191	0.9	<table border="1"><thead><tr><th>Product Name</th></tr></thead><tbody><tr><td>(Laptop) (HP191)</td></tr><tr><td>$t_1 \ 0.05 \times + 0.9 \times =$ (Attr. Sum. t_1)</td></tr></tbody></table>	Product Name	(Laptop) (HP191)	$t_1 \ 0.05 \times + 0.9 \times = $ (Attr. Sum. t_1)
Product Name																
Laptop HP191																
Word	Weight															
Laptop	0.05															
Notebook	0.07															
HP191	0.9															
Product Name																
(Laptop) (HP191)																
$t_1 \ 0.05 \times + 0.9 \times = $ (Attr. Sum. t_1)																
t_2	<table border="1"><thead><tr><th>Product Name</th></tr></thead><tbody><tr><td>Notebook HP191</td></tr></tbody></table>	Product Name	Notebook HP191		<table border="1"><thead><tr><th>Product Name</th></tr></thead><tbody><tr><td>(Notebook) (HP191)</td></tr><tr><td>$t_2 \ 0.07 \times + 0.9 \times =$ (Attr. Sum. t_2)</td></tr></tbody></table>	Product Name	(Notebook) (HP191)	$t_2 \ 0.07 \times + 0.9 \times = $ (Attr. Sum. t_2)								
Product Name																
Notebook HP191																
Product Name																
(Notebook) (HP191)																
$t_2 \ 0.07 \times + 0.9 \times = $ (Attr. Sum. t_2)																

Sequence-based attribute summarization

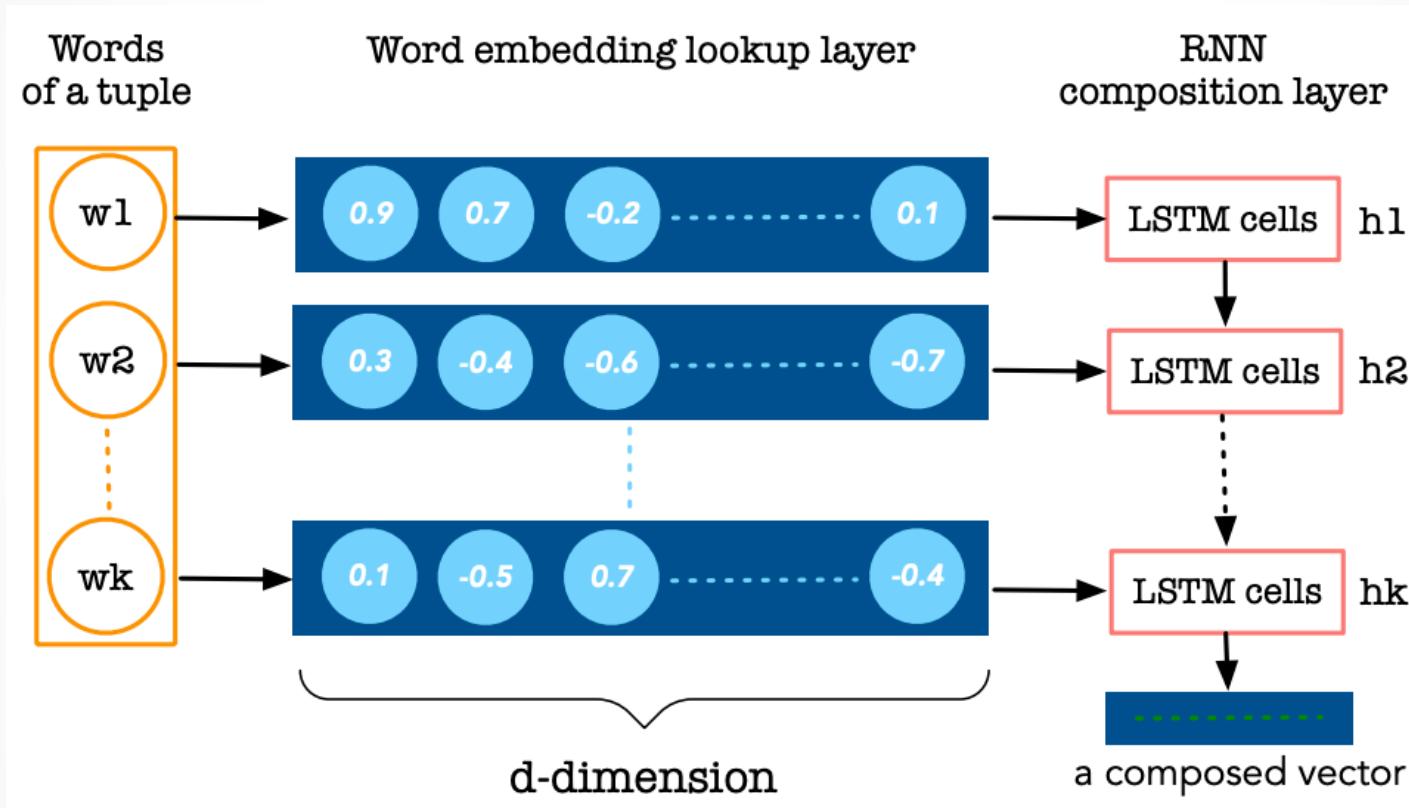
- **Vector captures which word sequences are present**
 - Automatically learns which word sequences are important

Product Name
black iPhone with gold case

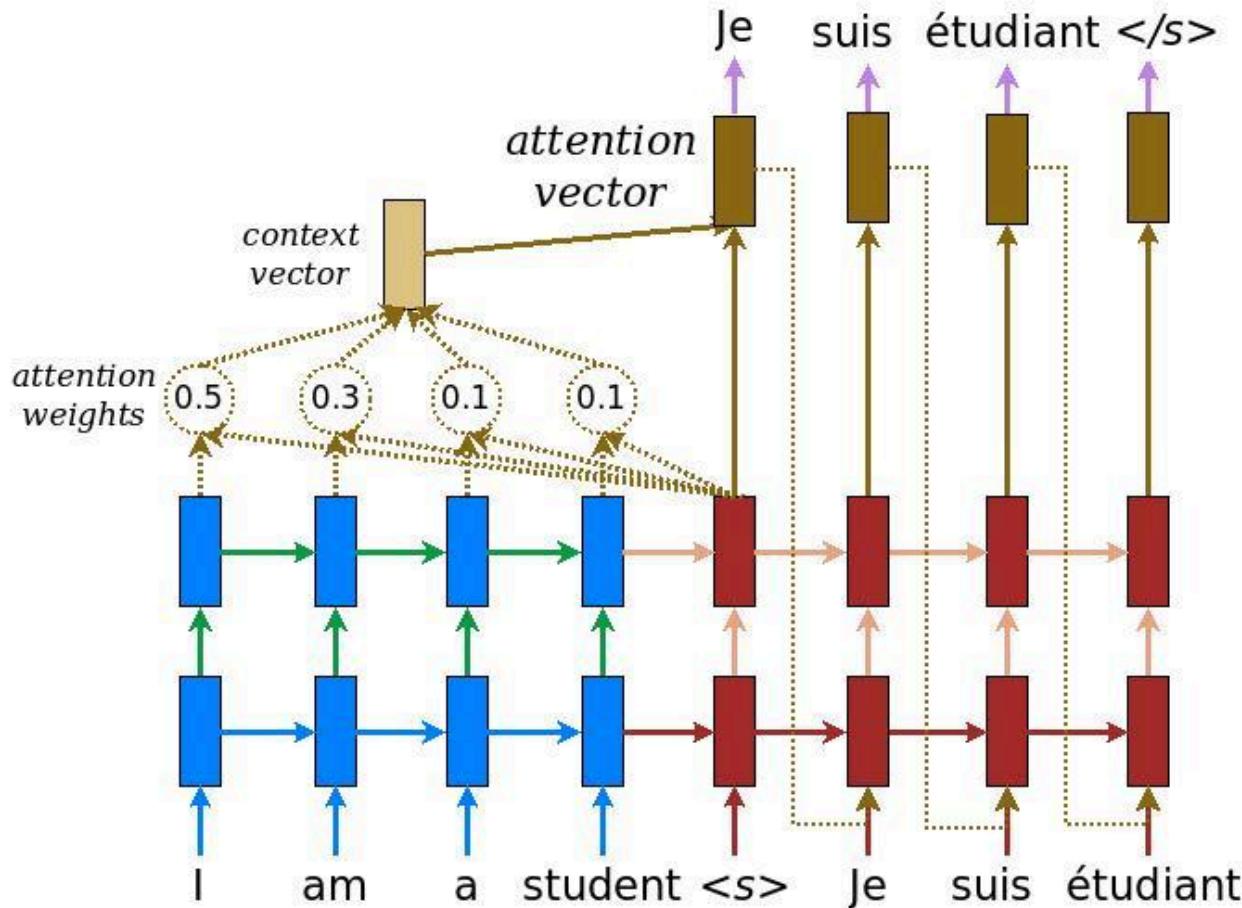
Product Name
gold iPhone with a black case



RNN: attribute summarization



Attention



Options at each stage

Architecture module		Options	
Attribute embedding		<i>Granularity:</i> (1) Word-based (2) Character-based	<i>Training:</i> (3) Pre-trained (4) Learned
Attribute similarity representation	(1) Attribute summarization	(1) Heuristic-based (2) RNN-based (3) Attention-based (4) Hybrid	
	(2) Attribute comparison	(1) Fixed distance (cosine, Euclidean) (2) Learnable distance (concatenation, element-wise absolute difference, element-wise multiplication)	
Classifier		NN (multi-layer perceptron)	

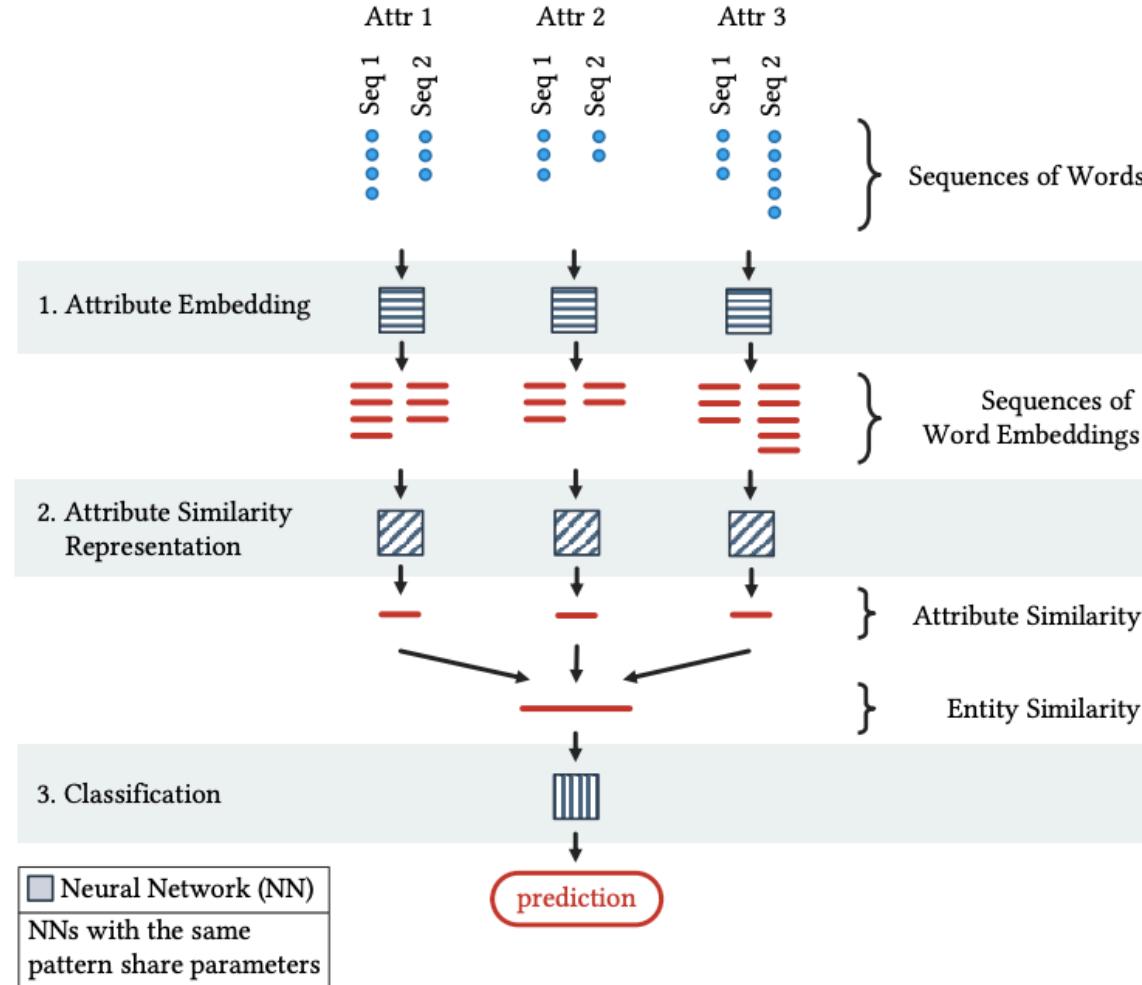
Figure 3: The design space of DL solutions for EM.

Results

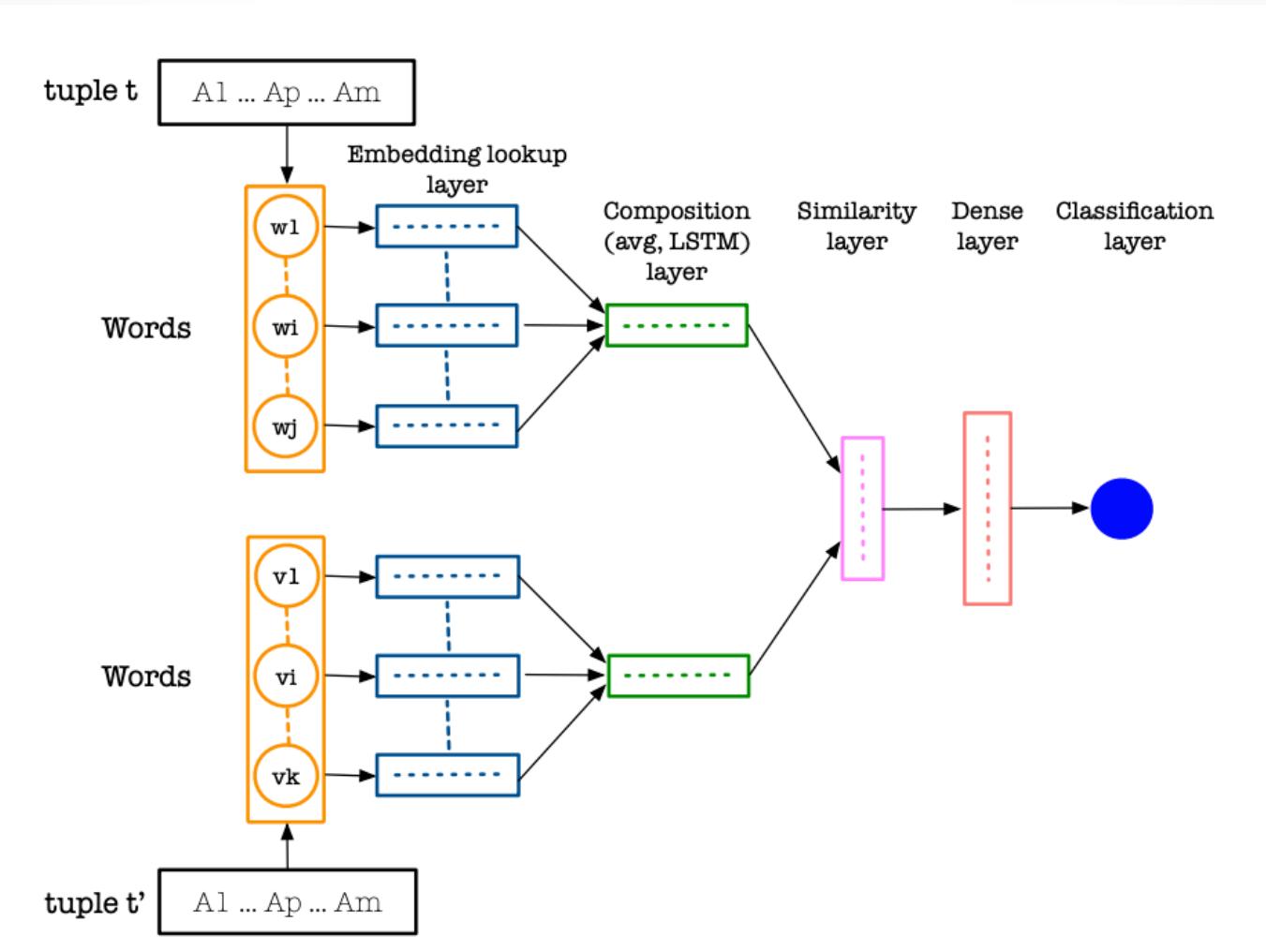
Table 9: F_1 -score comparison for different deep learning solutions and Magellan for different dataset types and sizes.

Dataset		SIF	RNN	Attention	Hybrid	Magellan
Structured	small	79.1	83.9	76.2	84.6	86.5
	large	91.0	91.8	91.8	91.5	91.3
Textual	small	35.1	39.4	56.8	62.8	43.6
	large	87.5	88.9	88.8	89.3	83.9
Textual w/o info. attr.	small	32.0	38.5	55.0	47.7	33.0
	large	86.1	86.9	88.4	89.7	80.2
Dirty	small	76.8	86.2	78.2	84.2	64.7
	large	85.6	89.6	90.3	90.0	72.3

Comparing Attributes



Comparing Entities



Results

Dataset	Magellan	DeepER	Published
Prod-WA	82.99	88.06	89.3 [26] (Crowd)
Prod-AG	87.68	96.029	62.2 [33] (ML)
Pub-DA	97.6	98.6	N/A
Pub-DS	98.84	97.67	92.1 [26] (Crowd)
Pub-DC	96.4	99.1	95.2 [17] (Crowd)
Rest-FZ	100	100	96.5 [26] (Crowd)