

# Extracting Meronyms for a Biology Knowledge Base Using Distant Supervision

Xiao Ling<sup>†</sup>

Peter Clark<sup>\*</sup>

Daniel S. Weld<sup>†</sup>

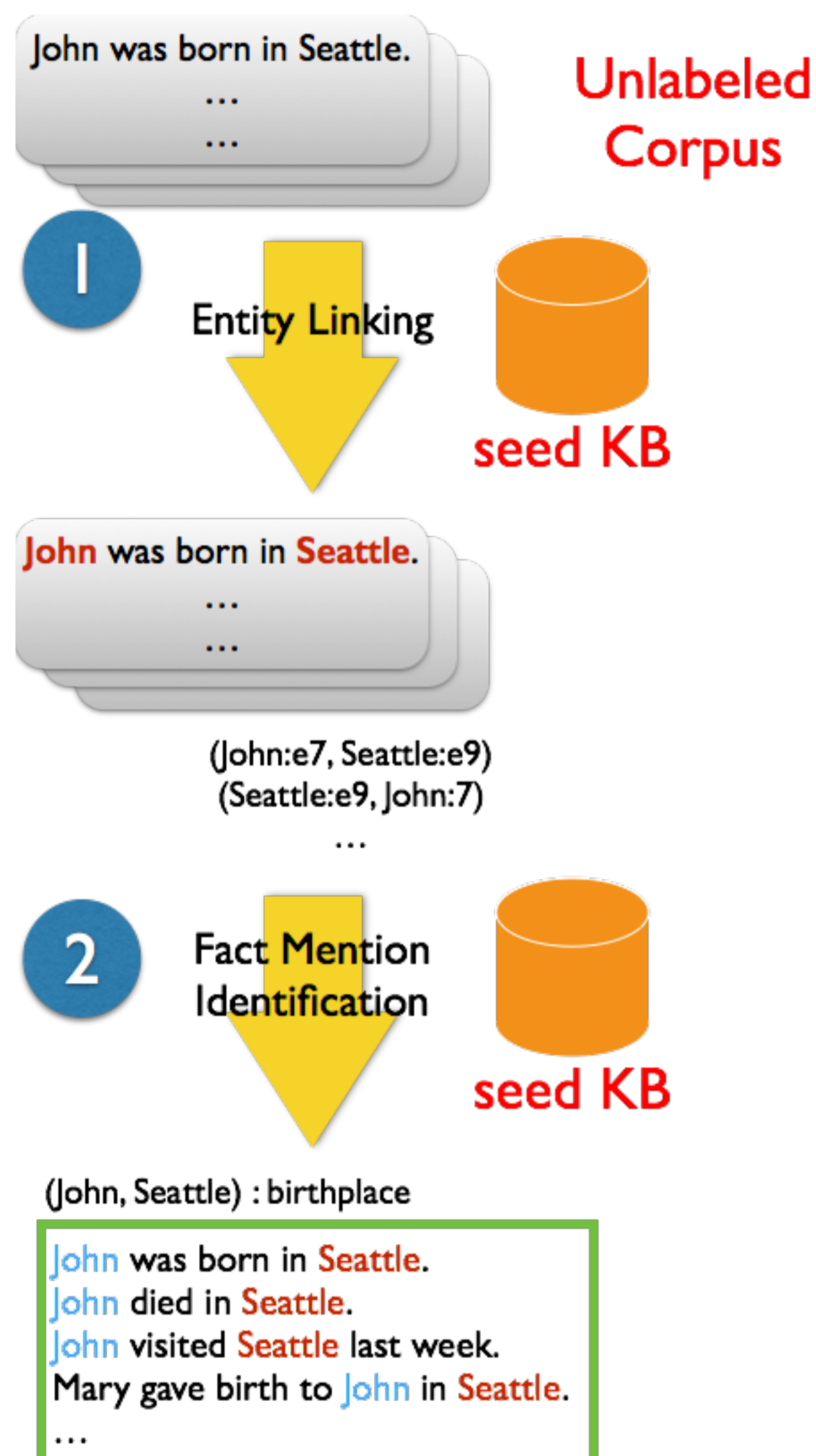
<sup>†</sup>University of Washington

<sup>\*</sup>Allen Institute for Artificial Intelligence

## Overview

- ▶ Goal: Build a **KB** of *meronym* relations for the domain of *biology*
- ▶ Methodology: use **distant supervision** to learn a meronym relation extractor.
- ▶ Data: 1) a seed KB of meronym facts; 2) a unlabeled text corpus (i.e. a textbook “Campbell Biology”).
- ▶ Evaluation: a held-out set (172 has-part, 206 NA)

## Distant Supervision



## Expanding mentions with co-reference

- ▶ Standard: NER + string match
- ▶ Named Entity Linking
- ▶ This work: dictionary max-span matching

	Recall	Precision	F1
CV	0.664	0.820	0.733
CV+COREF	<b>0.674</b>	<b>0.821</b>	<b>0.740</b>
TEST	0.663	<b>0.857</b>	0.748
TEST+COREF	<b>0.744</b>	0.795	<b>0.769</b>

Example:

One of the sex-determining chromosomes is **X chromosome**; and it has around 1,100 **genes**.

X-Chromosome → Gene

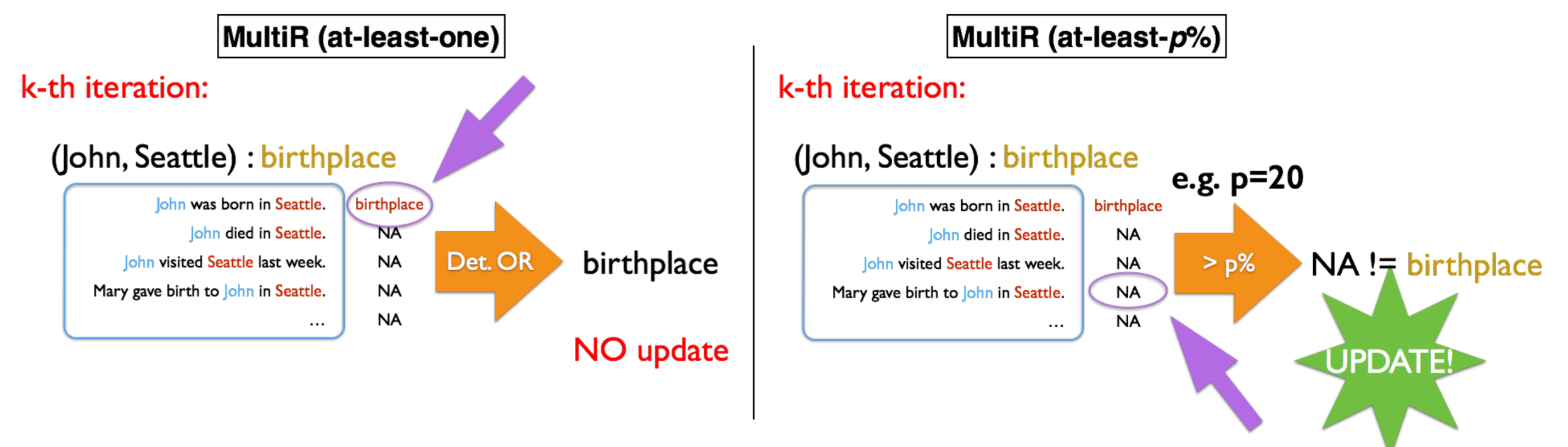
## Generating negative data

- ▶ The seed KB provides mostly valid facts while invalid facts can be useful as negative data.
- ▶ To generate negative examples:
  - ▷ random sample (closed world)
  - ▷ functional relations

	#+	#-	Recall	Precision	F1
CV(BASE)	887	77	<b>0.856</b>	0.389	0.533
CV(REV)	887	963	0.706	0.748	0.725
CV(TRANS)	887	2566	0.674	<b>0.821</b>	<b>0.740</b>
TEST(BASE)	887	77	<b>0.884</b>	0.596	0.712
TEST(REV)	887	963	0.709	0.718	0.713
TEST(TRANS)	887	2566	0.744	<b>0.795</b>	<b>0.769</b>

- This work:
- ▶ **reverse:** if  $(e_1, \text{has-part}, e_2)$  then  $(e_2, \text{NA}, e_1)$
  - ▶ **transitive closure**

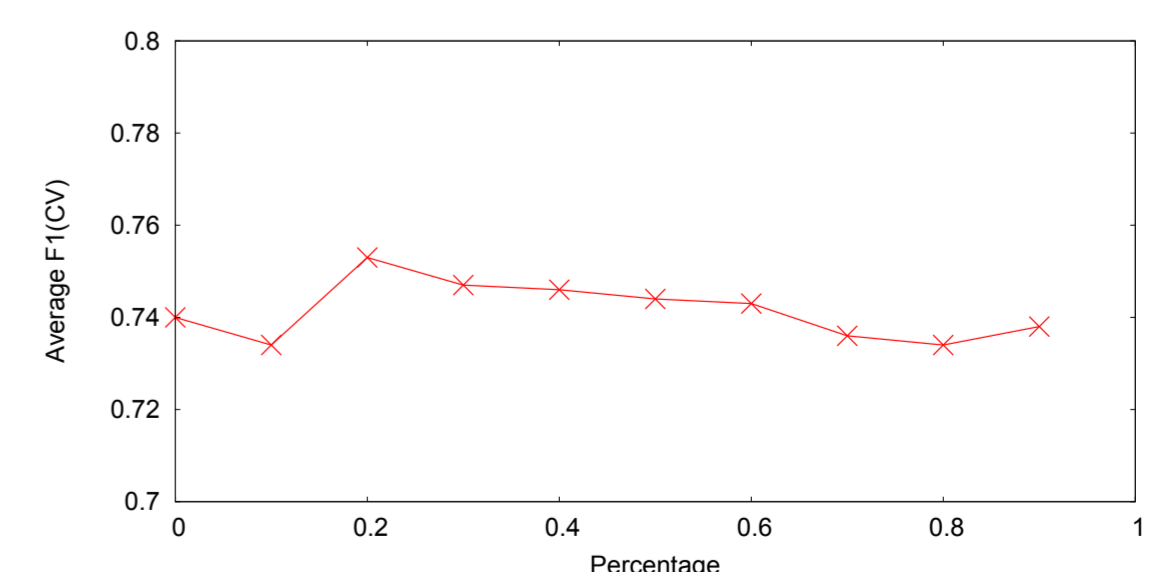
## at-least- $p\%$ assumption



## Leveraging supervision from out-of-domain

- ▶ additional training facts: meronyms from WordNet
- ▶ We match the WordNet meronyms to Wikipedia articles.

	p	#+	#-	Recall	Prec	F1
TEST	0.0	887	2566	0.384	<b>0.892</b>	0.537
TEST+WN	0.0	1578	2566	<b>0.401</b>	0.841	<b>0.543</b>
TEST	0.2	887	2566	0.523	<b>0.811</b>	0.636
TEST+WN	0.2	1578	2566	<b>0.558</b>	0.793	<b>0.655</b>



	Recall	Precision	F1
CV( $p = 0$ )	0.674	<b>0.821</b>	0.740
CV( $p = 0.2$ )	<b>0.730</b>	0.776	<b>0.753</b>
TEST( $p = 0$ )	0.744	<b>0.795</b>	0.769
TEST( $p = 0.2$ )	<b>0.791</b>	0.786	<b>0.788</b>