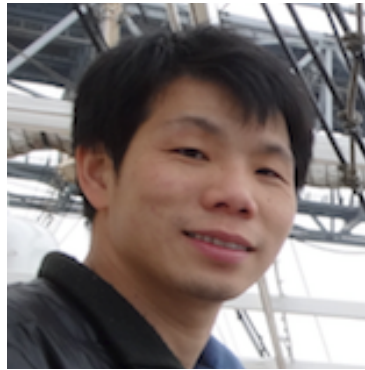


Cross-lingual Open Information Extraction

Sheng Zhang



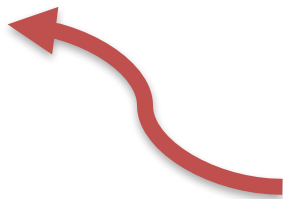
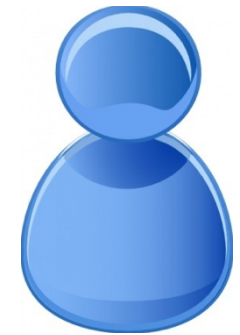
Kevin Duh



Ben Van Durme

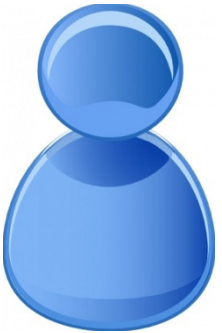
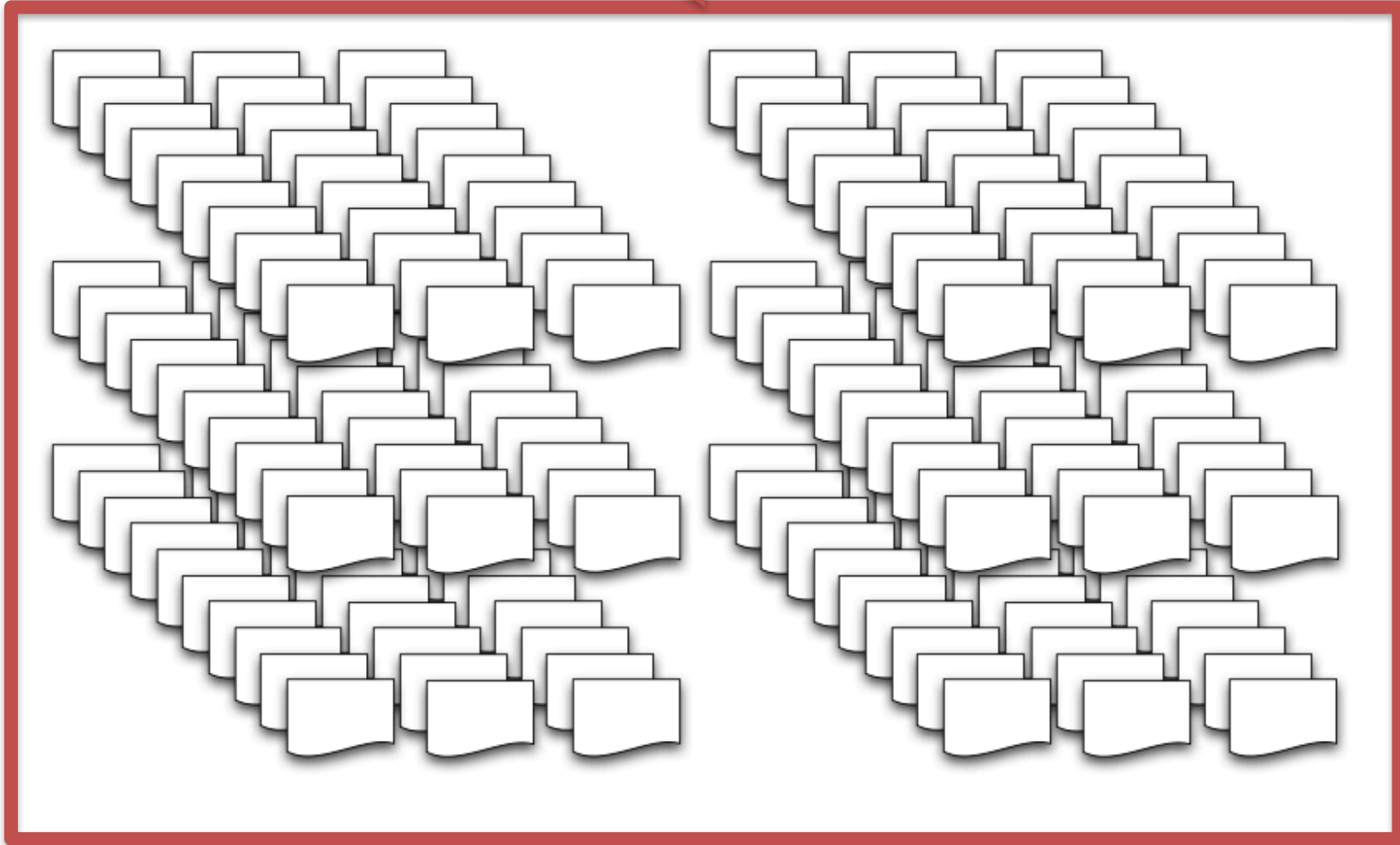


Johns Hopkins University

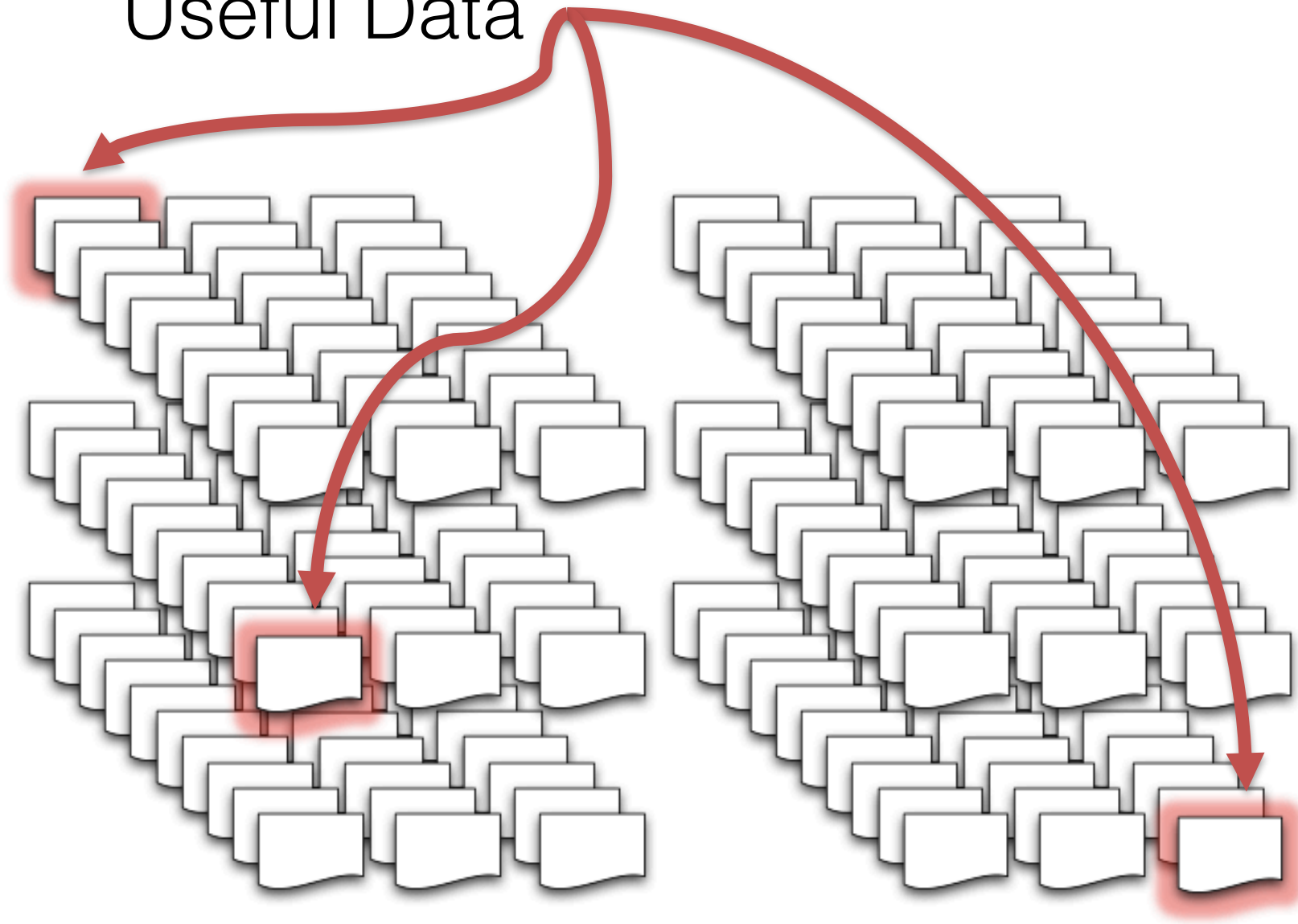


Analyst: User with complex information need

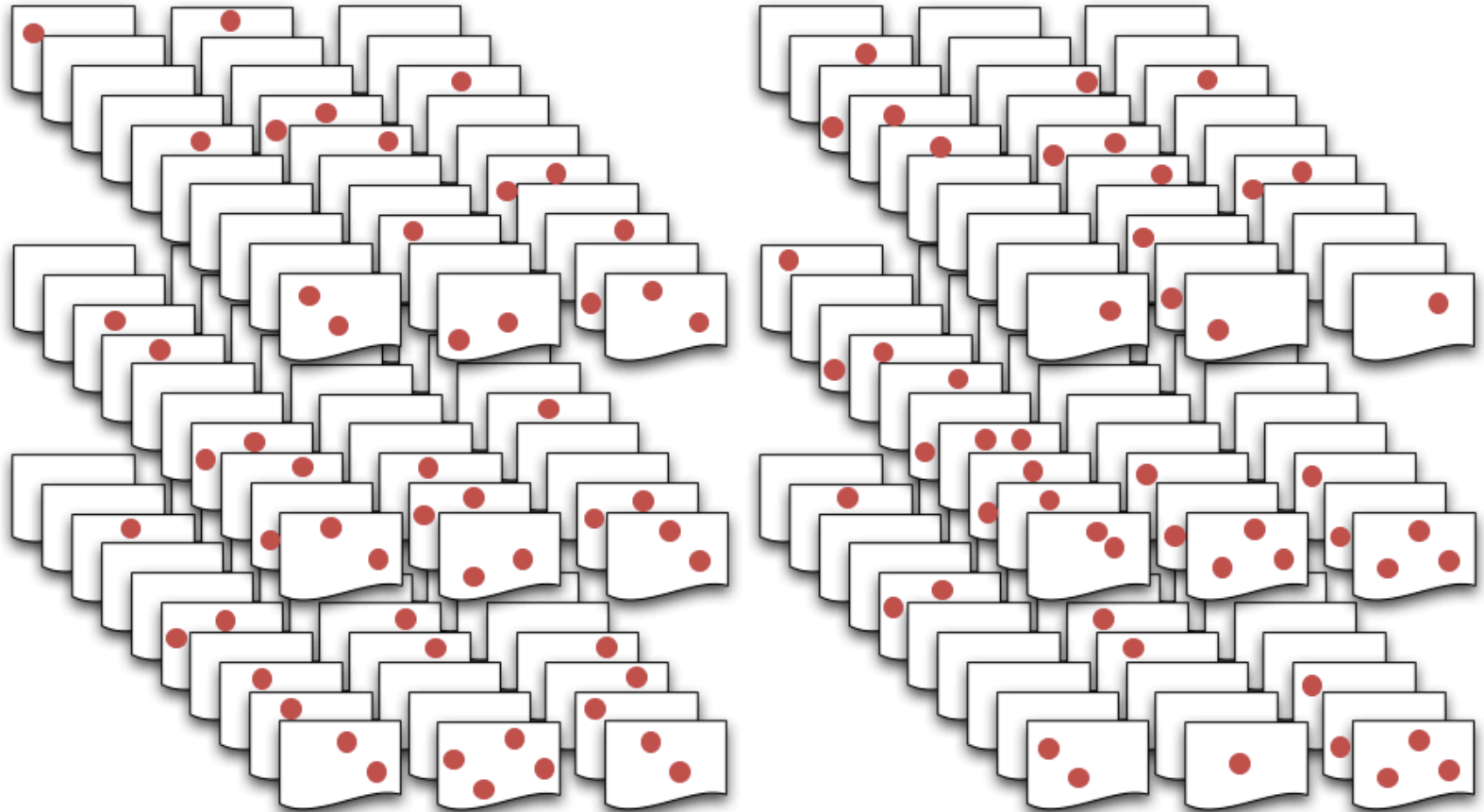
Data



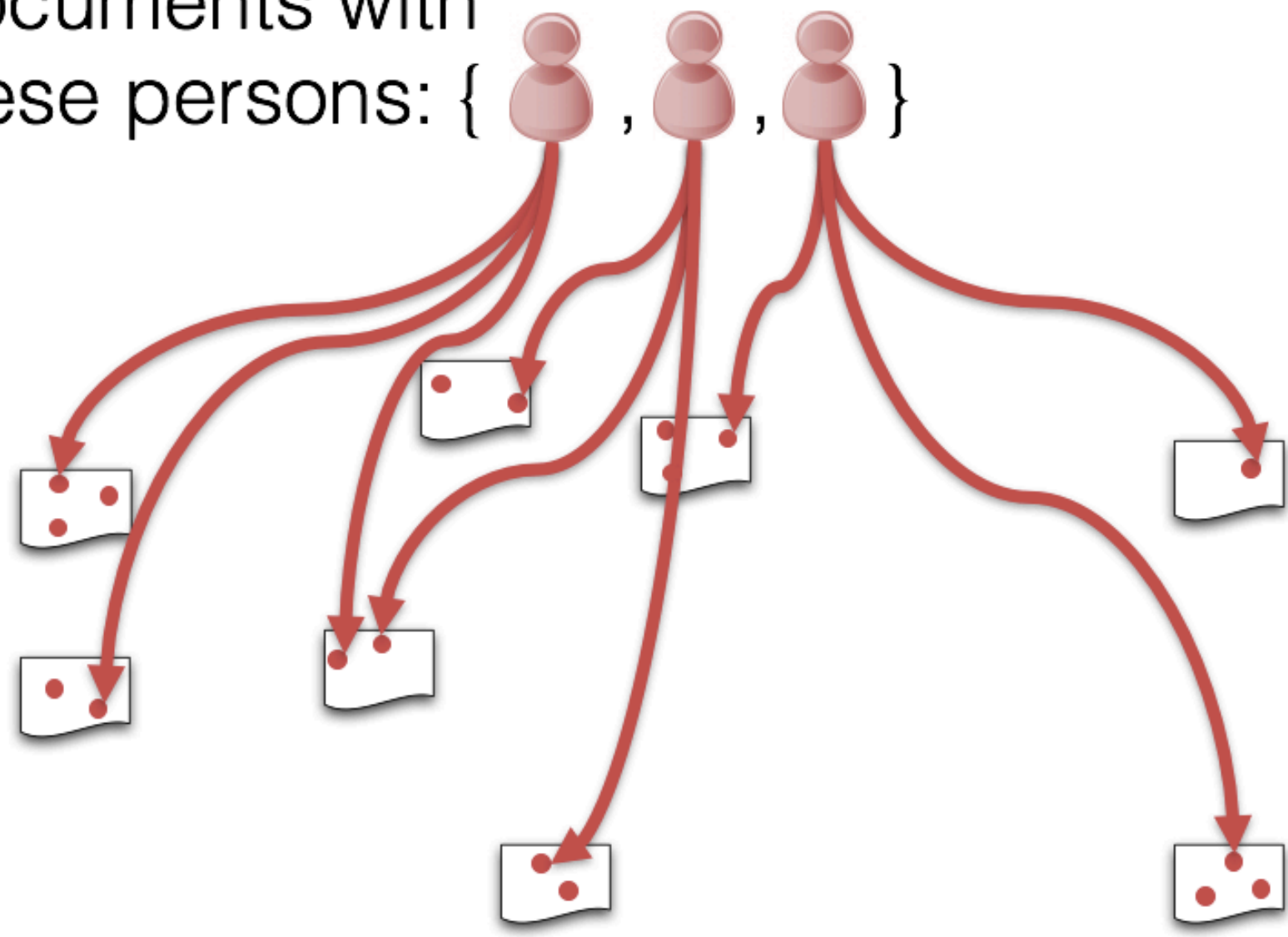
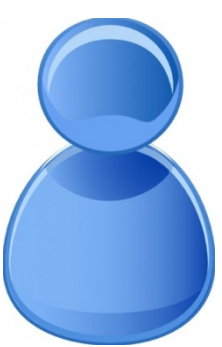
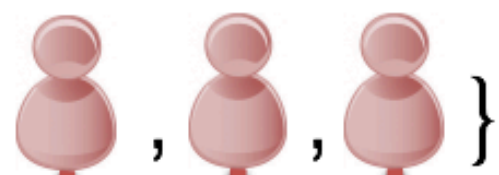
Useful Data



Entity Recognition

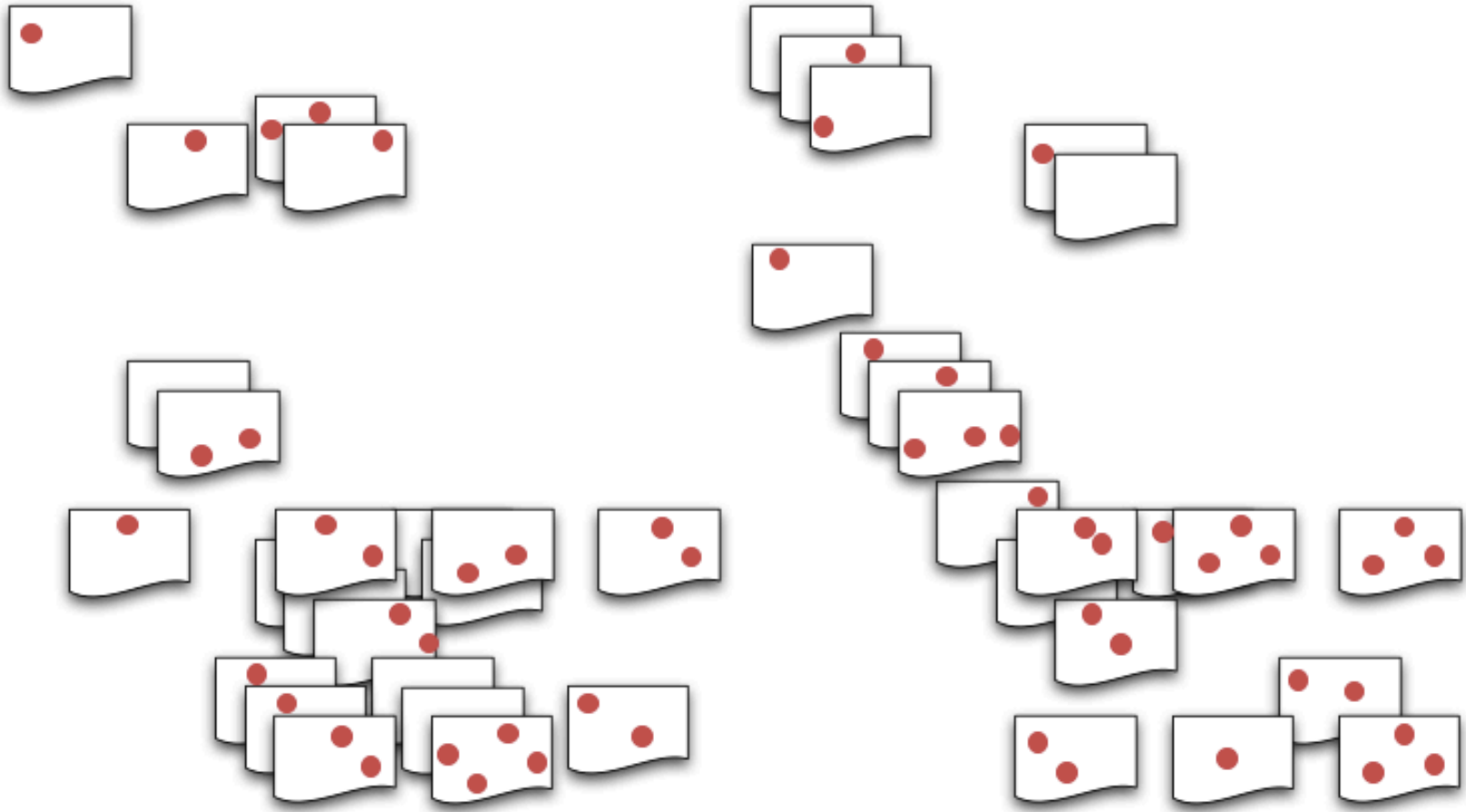


Documents with these persons: {



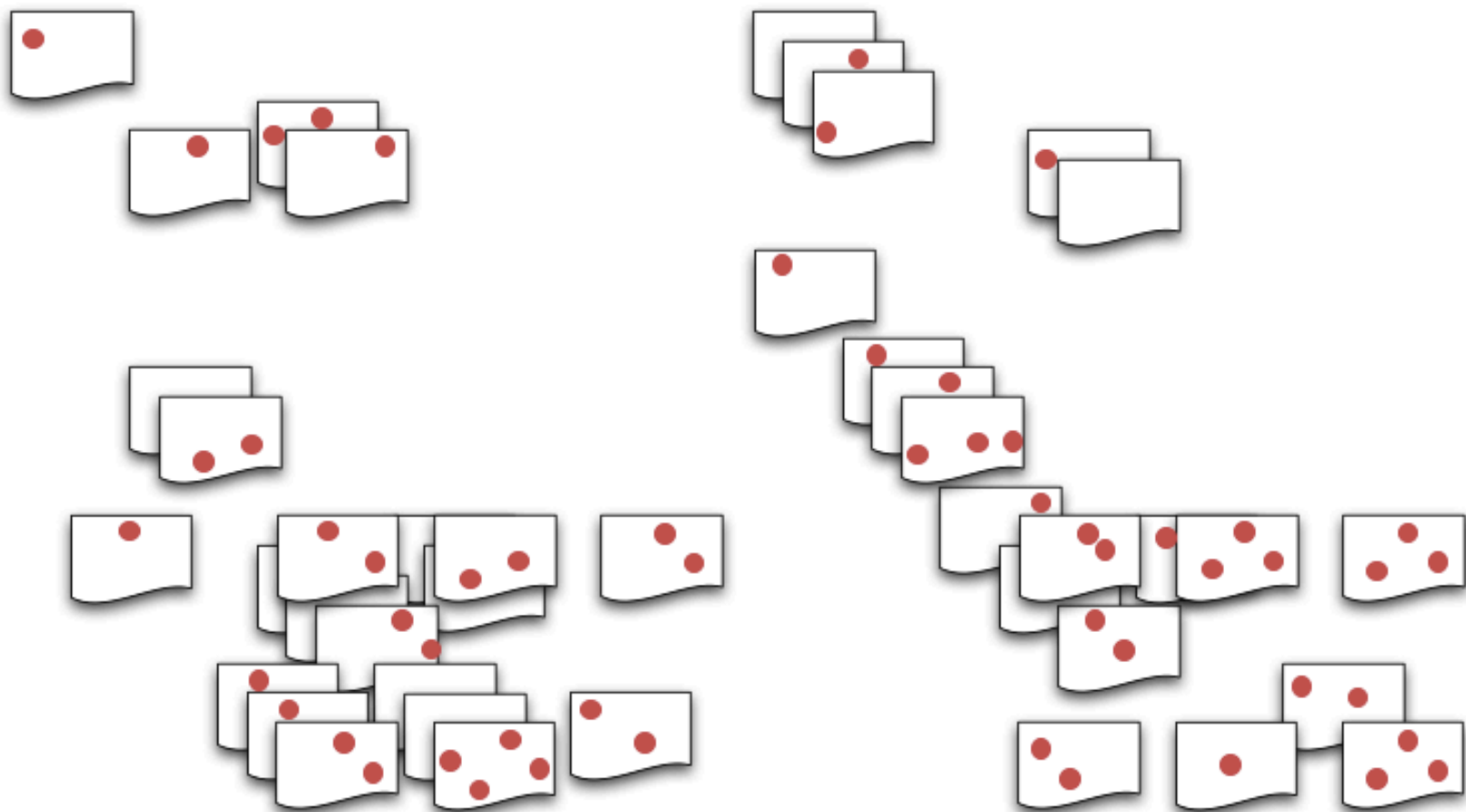
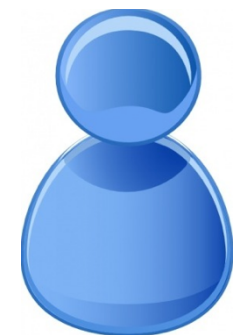
Documents with a:

CausalRelation(DISEASE, SYMPTOM)



Documents with a:

CausalRelation(DISEASE, SYMPTOM)



Foreign Language Documents

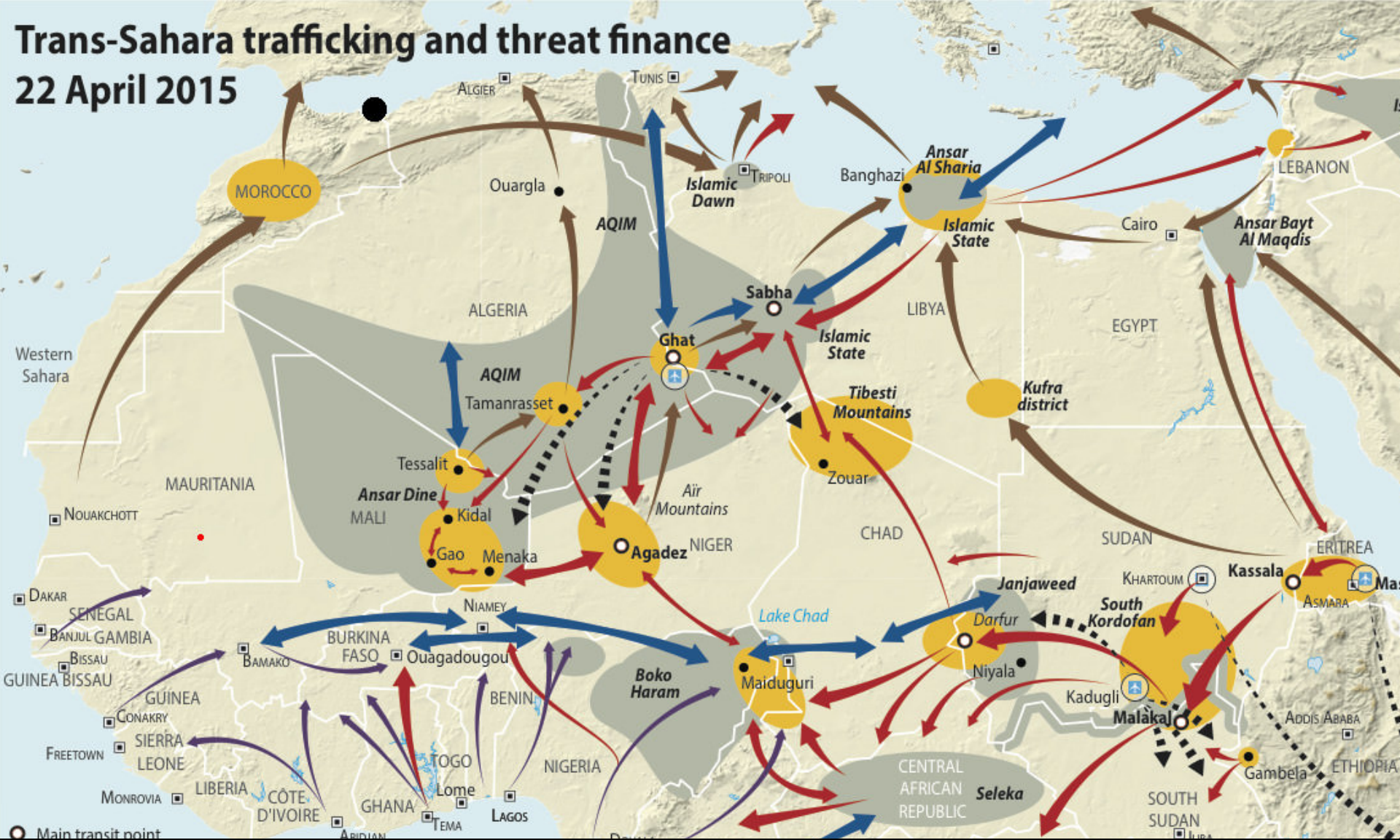


Financial Analyst: Do I invest in commodity futures?



Trans-Saharan trafficking and threat finance

22 April 2015



Intelligence Analyst: How are the terrorists connected?





Aid Worker: Which locations have immediate need?



Outline

1. Motivation
2. Problem Definition
3. Pipeline vs. Joint Solution
4. Improved Joint Solution

Task Formulation

- Cross-lingual:
 - analyst speaks English, but document collection is in other languages
- Cross-lingual **Information Retrieval?**
 - Document unit is too large
- Cross-lingual **Question Answering?**
 - Difficulty in formulating questions
- Cross-lingual **Information Extraction?**
 - Close, but no fixed ontology

Information Extraction vs. **Open** Information Extraction

Bill Gates, Microsoft co-founder, stepped down as CEO in January 2000. **Gates** was included in the **Forbes wealthiest list** since 1987 and **was the wealthiest** from 1995 to 2007...

It was announced that **IBM** would buy **Ciao** for an undisclosed amount. The **CEO, MacLorraine** has occupied the **corner office of the Hopkinton**, company

The company's storage business is also threatened by new, born-on-the Web could providers like Dropbox and Box, and ...

IE
→

Co-founder(Bill Gates, Microsoft)
Director-of (MacLorraine, Ciao)
Employee-of (MacLorraine, Ciao)
...

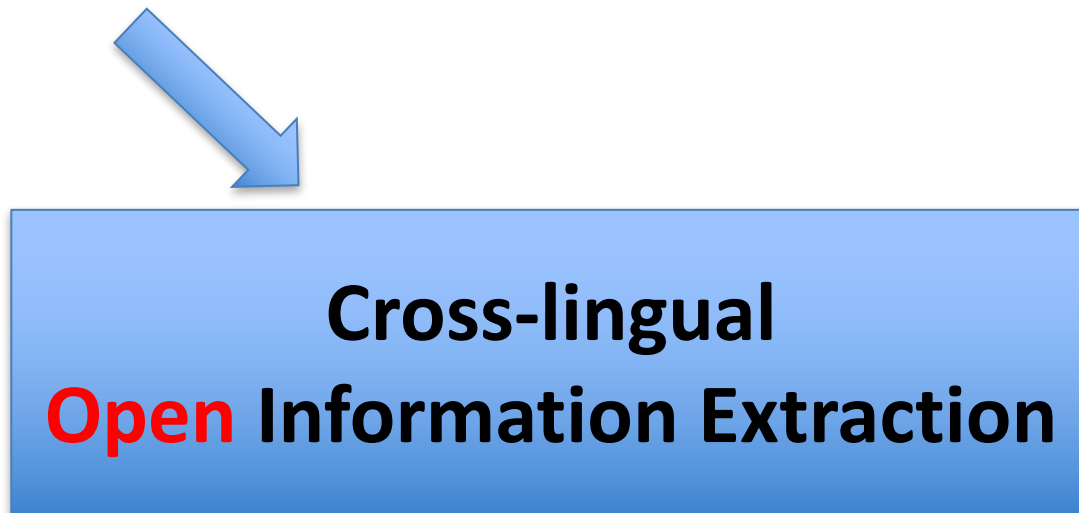
Open IE
→

(Bill Gate, be, Microsoft co-founder)
(Bill Gates, stepped down as, CEO)
(Bill Gates, was included in, the Forbes wealthiest list)
(Bill Gates, was, the wealthiest)
(IBM, would buy, Ciao)
(MacLorraine, has occupied, the corner office of the Hopkinton)
...

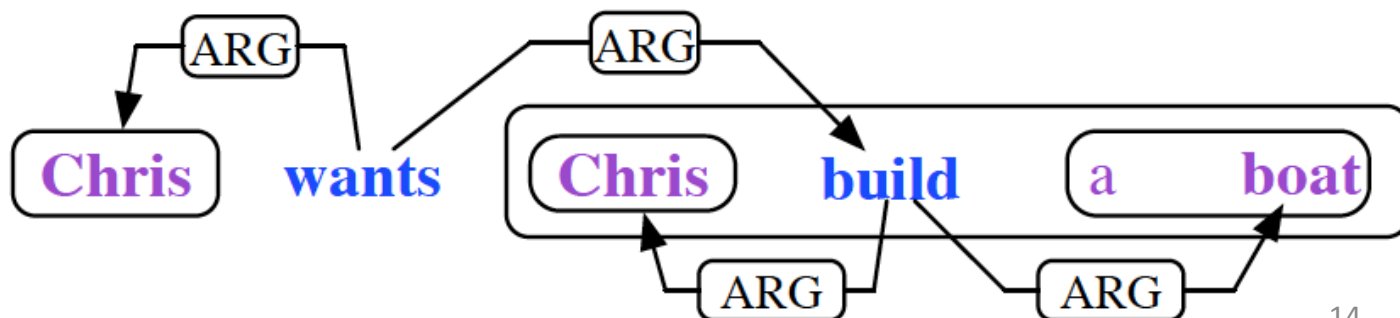
	IE	Open IE
Input	Sentences + Labeled relations	Sentences
Relation	Specified relations in advance	Free discovery
Extractor	Specified relations	Independent-relations

Input: Chinese sentence

克里斯想造一艘船。



Output: A set of English tuples, e.g. Relation(arg1,arg2)



Cross-lingual **Open** Information Extraction

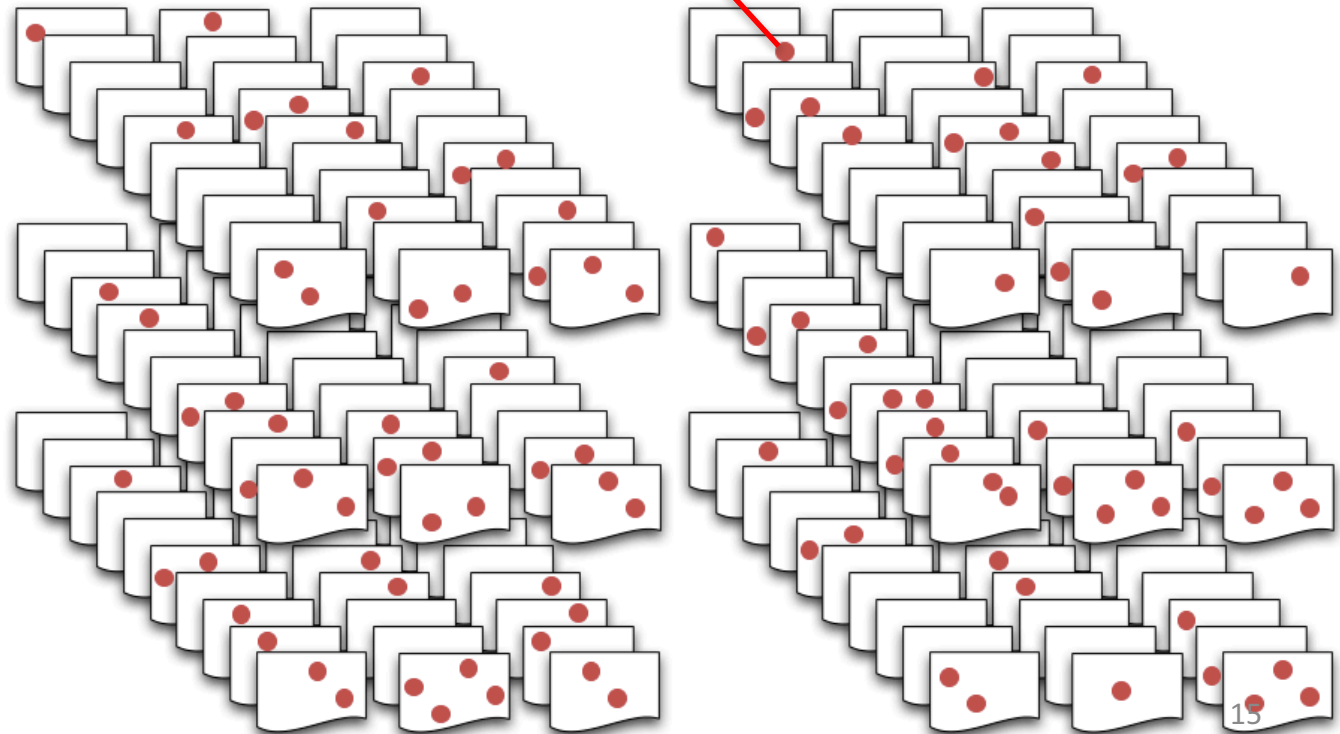
Visualization

RelationA(arg1,arg2)

RelationB(arg1,arg2)

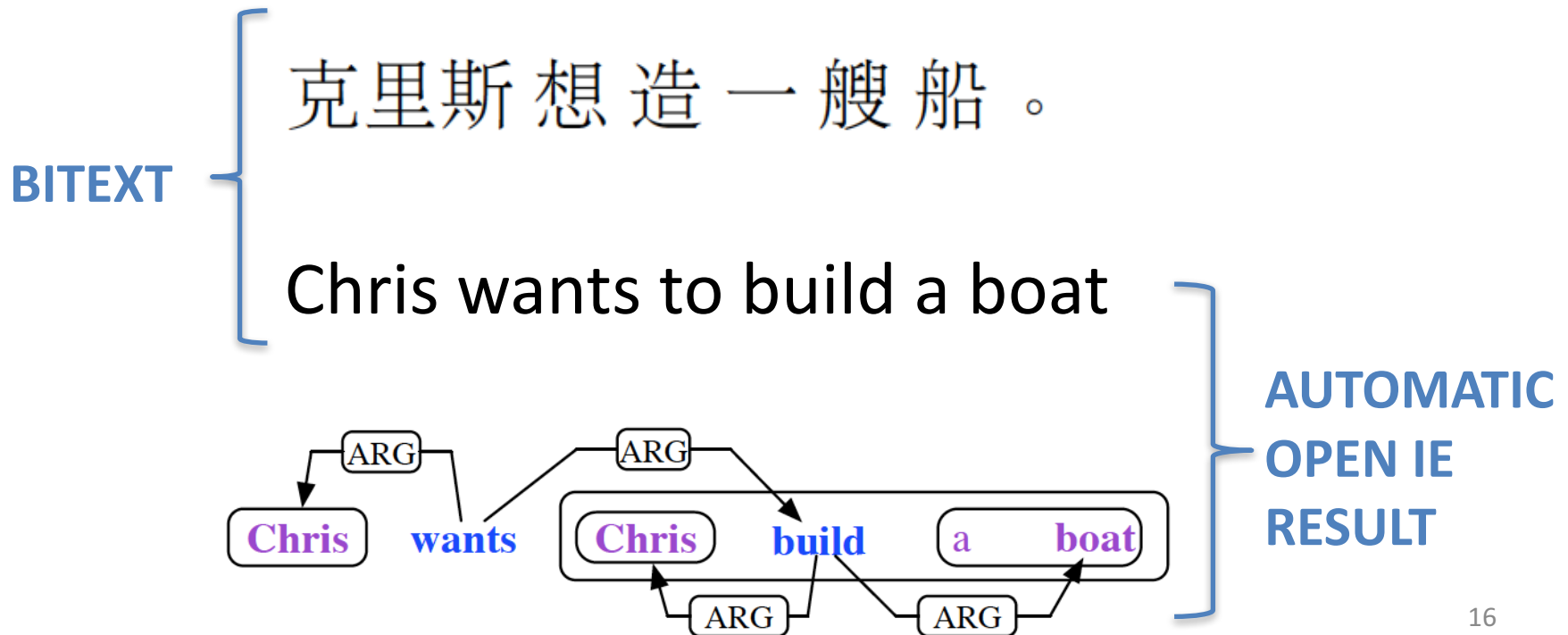
RelationC(arg1,arg2)

Query



Assumptions

1. Training data: Chinese-English bitext
2. Monolingual Open IE system in English



Monolingual Open IE System

PredPatt: <https://github.com/hltcoe/PredPatt>

- Based on Universal Dependencies
- Rules for:
 1. identifying **predicate root** and **argument root**:
e.g. nsubj(**s**, **v**), dobj(**o**, **v**)
 2. resolving arguments:
Chris expects to visit Pat → nsubj(Chris,visit)
Chris likes to sing and dance → nsubj(Chris,dance)
 3. phrase extraction:
PredPatt finds structure in text → ?a finds ?b in ?c

Pierre Vinken , 61 years old , will join the board as a nonexecutive director Nov. 29 .

?a is/are 61 years old

?a: Pierre Vinken

?a will join ?b as ?c ?d

?a: Pierre Vinken , 61 years old

?b: the board

?c: a nonexecutive director

?d: Nov. 29

?a is/are nonexecutive

?a: a director

Mr. Vinken is chairman of Elsevier N.V. , the Dutch publishing group .

?a is chairman of ?b

?a: Mr. Vinken

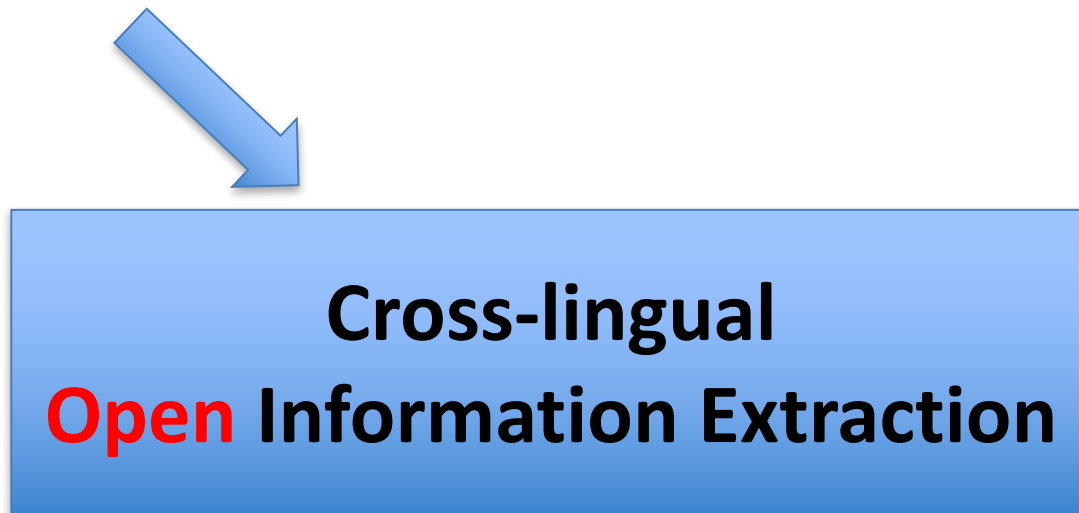
?b: Elsevier N.V.

?a is/are the Dutch publishing group

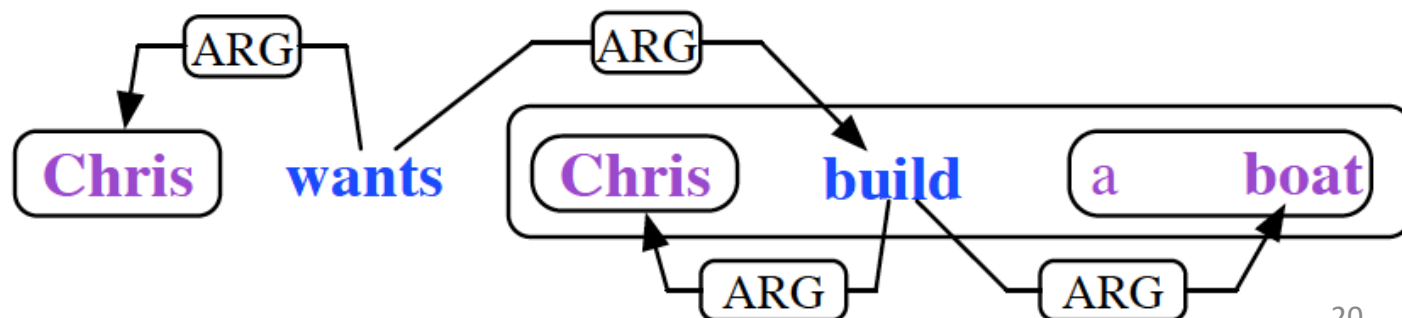
?a: Elsevier N.V.

Input: Chinese sentence

克里斯想造一艘船。



Output: A set of English tuples, e.g. Relation(arg1,arg2)

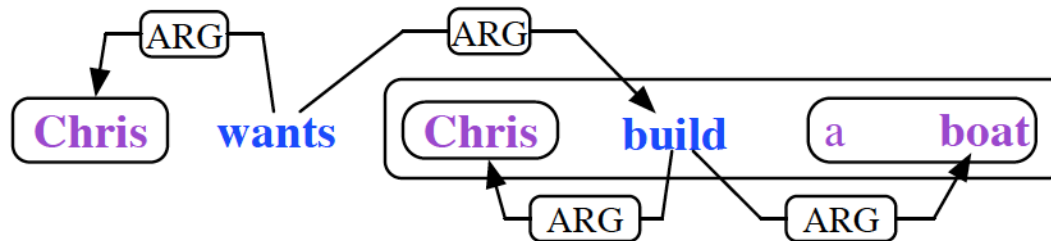


Outline

1. Motivation
2. Problem Definition
3. Pipeline vs. Joint Solution
4. Improved Joint Solution

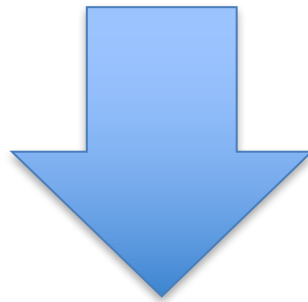
克里斯想造一艘船。

Chris wants to build a boat



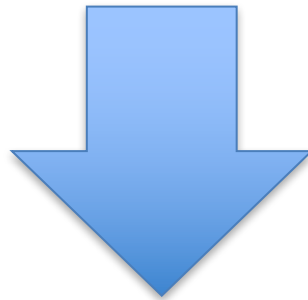
PIPELINE SOLUTION

克里斯想造一艘船。

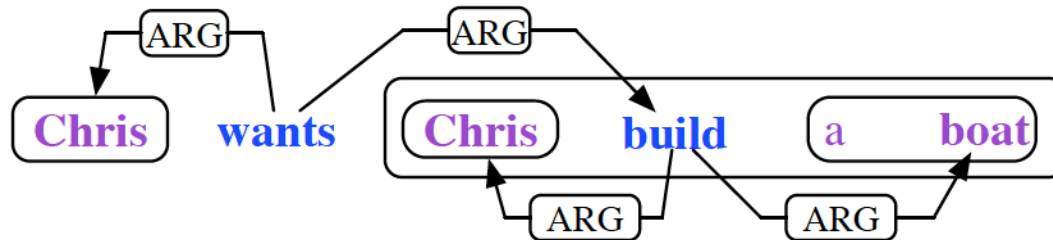


Machine Translation

Chris wants to build a boat



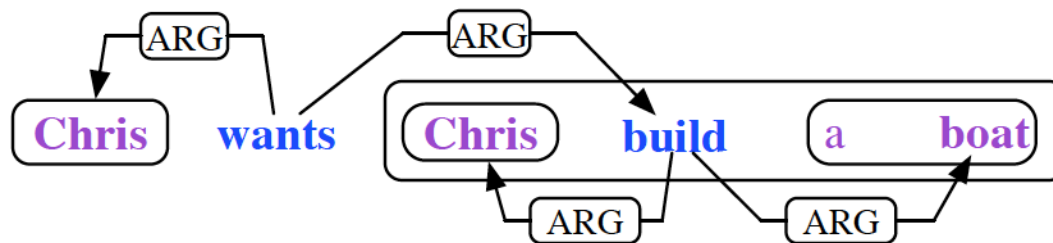
Dependency Parser +
English Open IE



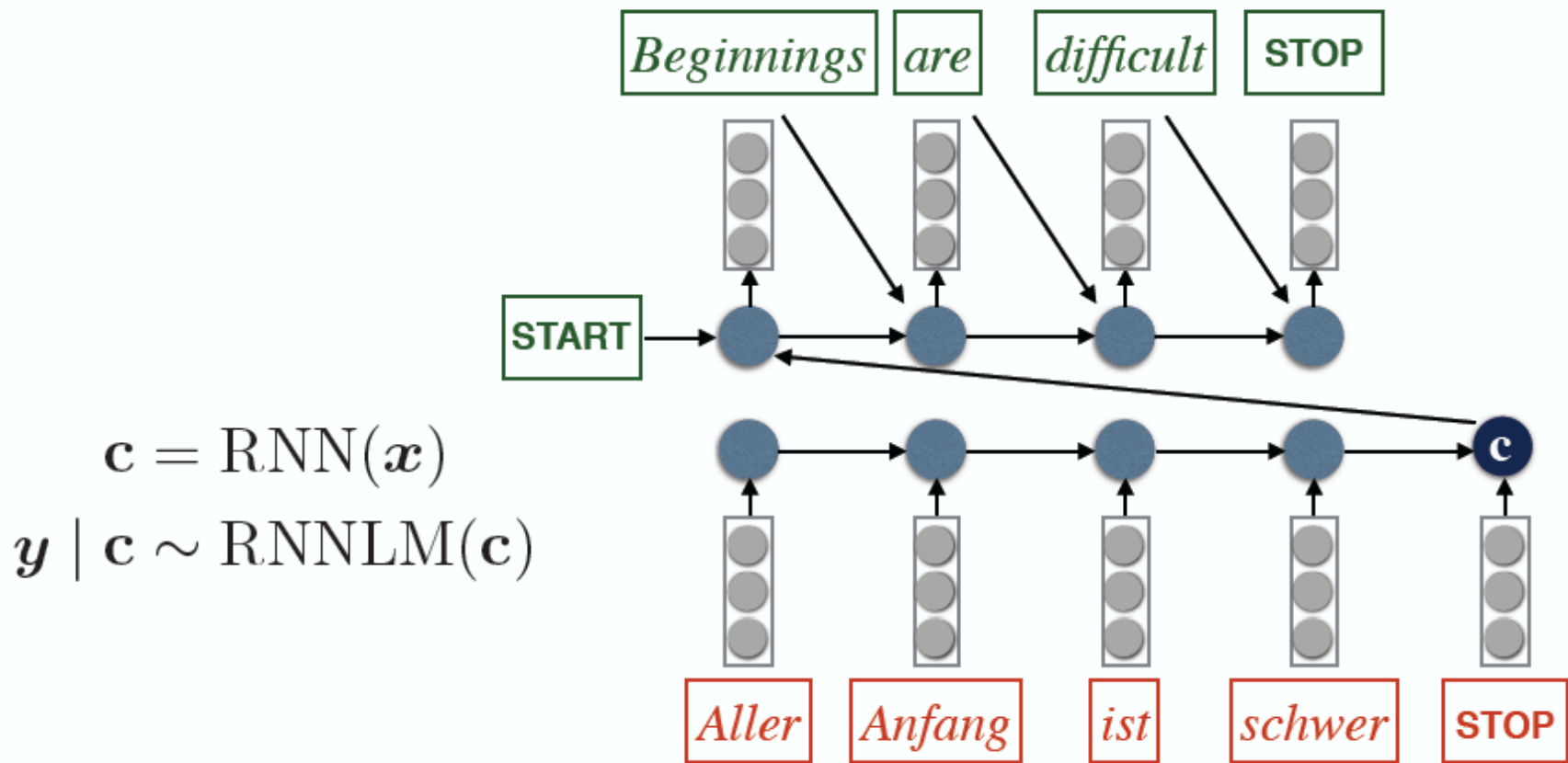
JOINT SOLUTION

克里斯想造一艘船。

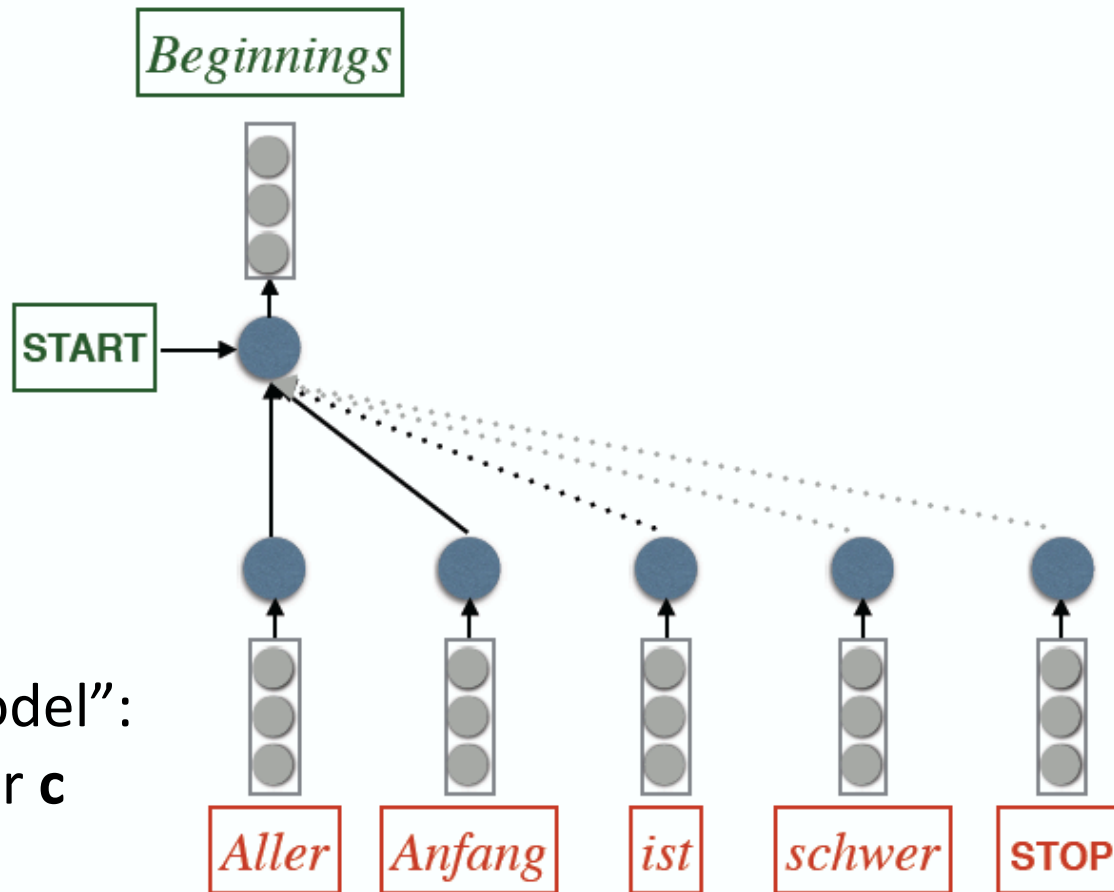
Cross-lingual Open IE



Neural Sequence-to-Sequence Model

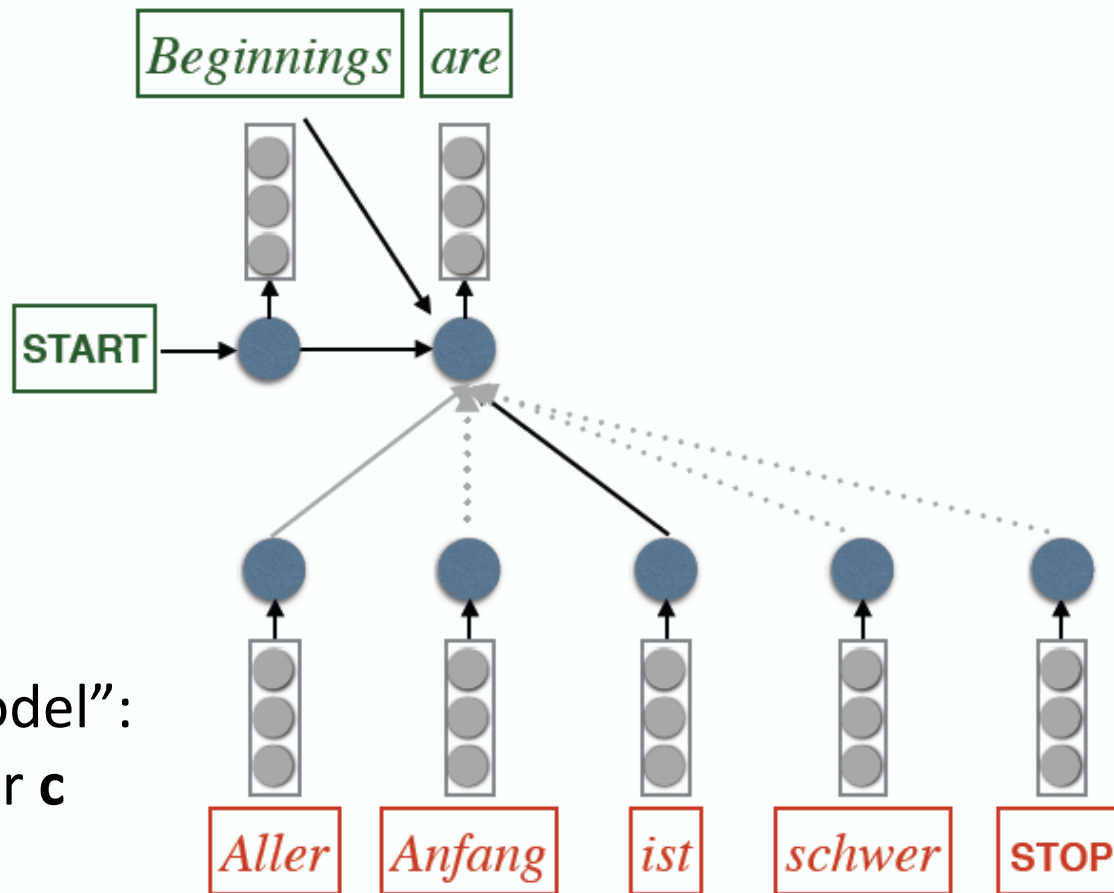


Neural Sequence-to-Sequence Model



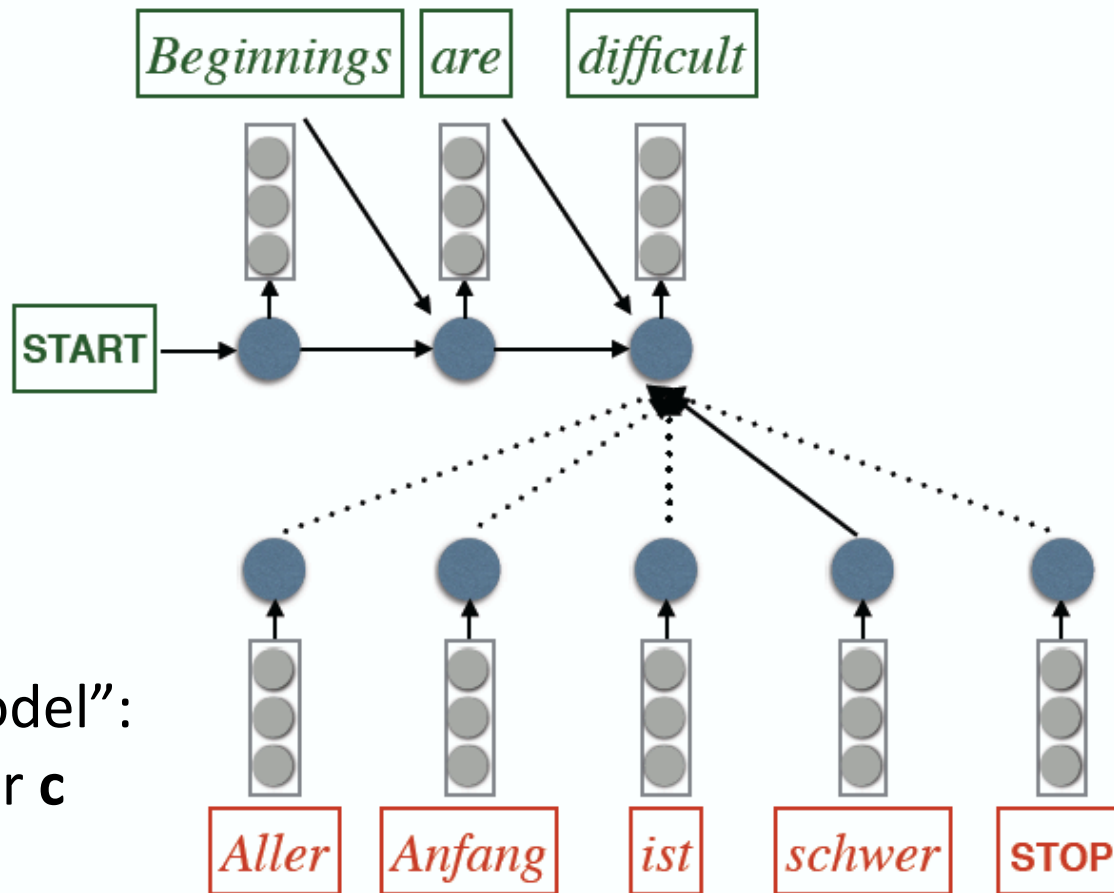
“Attention model”:
Context vector \mathbf{c}
is dynamic

Neural Sequence-to-Sequence Model



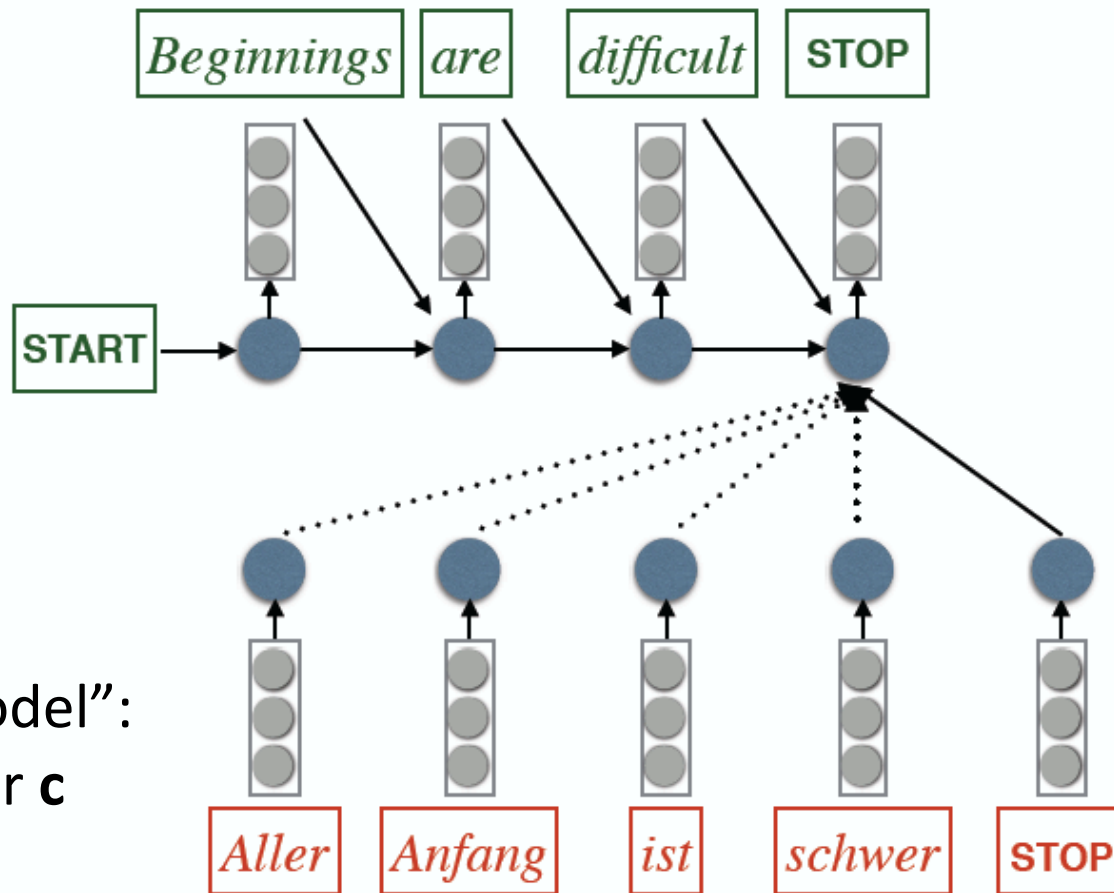
“Attention model”:
Context vector \mathbf{c}
is dynamic

Neural Sequence-to-Sequence Model



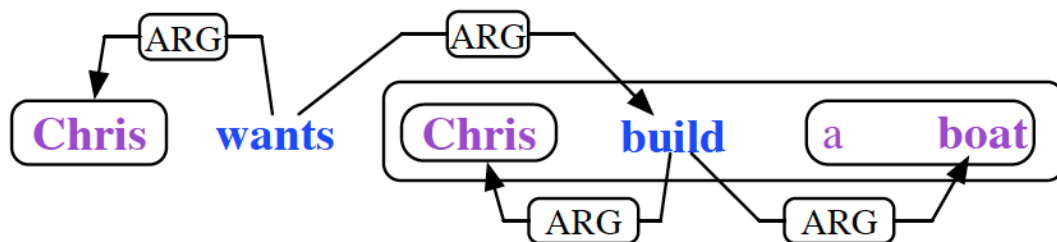
“Attention model”:
Context vector \mathbf{c}
is dynamic

Neural Sequence-to-Sequence Model

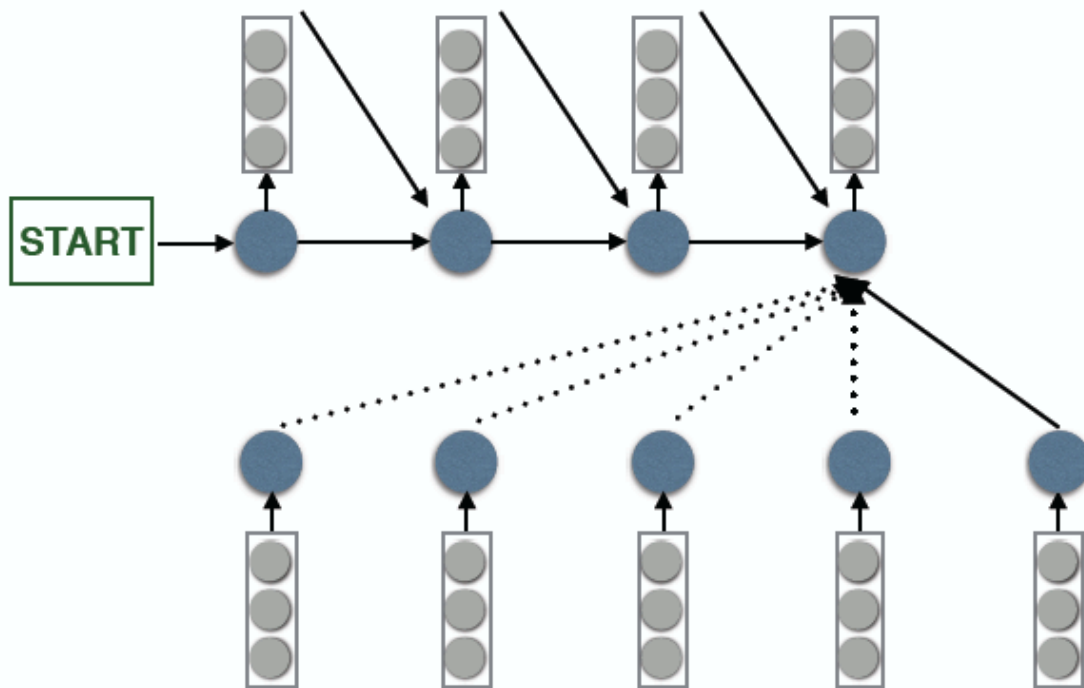


“Attention model”:
Context vector \mathbf{c}
is dynamic

Linearized OpenIE output as target



$[(\text{Chris}:a_h) \text{ wants}:p_h [(\text{Chris}:a_h) \text{ build}:p_h (\text{a}:a \text{ boat}:a_h)]]$

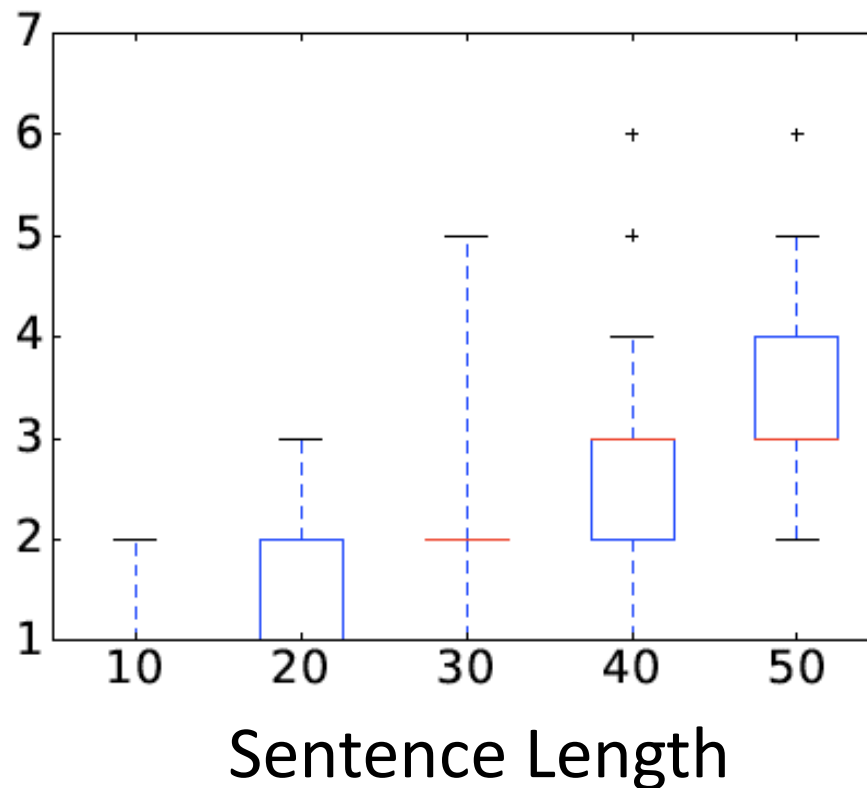


克里斯想造一艘船。

Experiment Setting

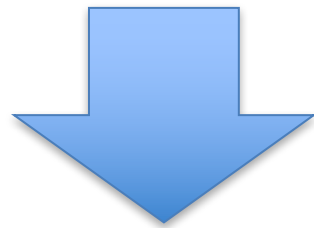
- 1 million sentence Chinese-English bitext (GALE project; mixed domain)

Predicates
in OpenIE
output



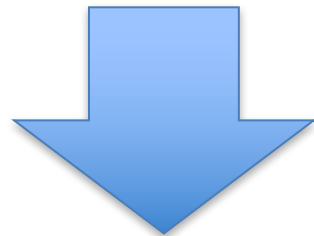
PIPELINE: BLEU=17.2 / PredicateF1=24.2

克里斯想造一艘船。

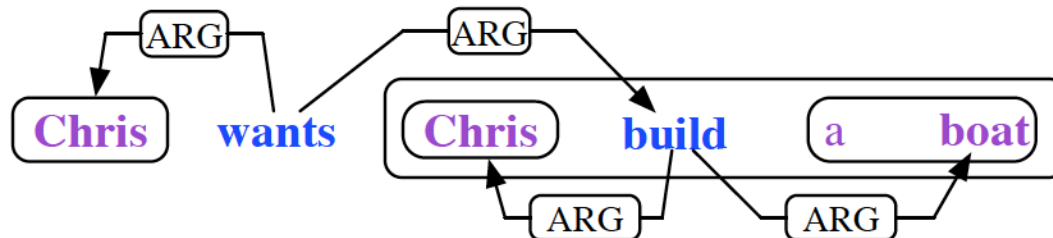


Machine Translation (Moses)

Chris wants to build a boat



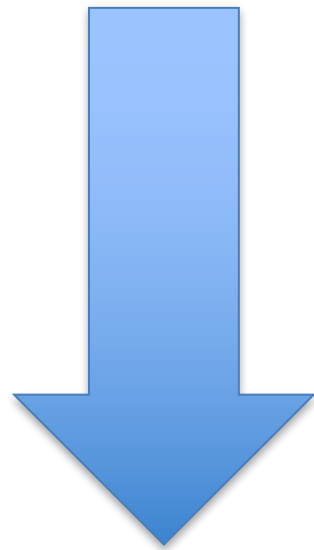
Dependency Parser (Parsey) +
English Open IE (PredPatt)



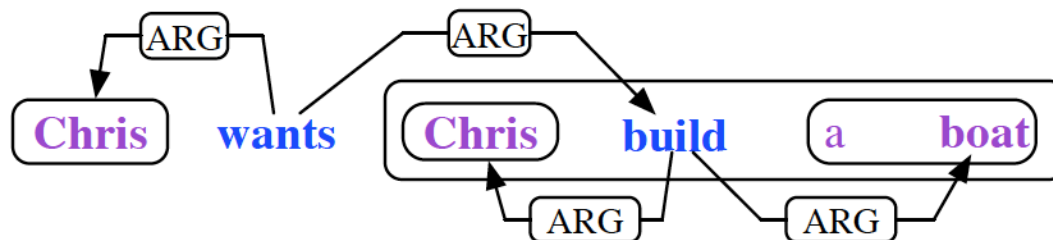
PIPELINE: BLEU=17.2 / PredicateF1=24.2

JOINT w/ Moses: BLEU=18.3 / PredicateF1=25.1

克里斯想造一艘船。



**Phrase-based
Machine Translation
(Moses)**

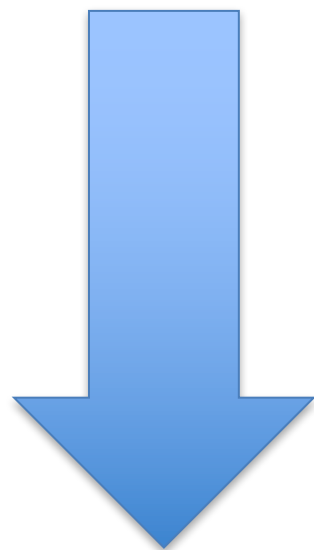


PIPELINE: BLEU=17.2 / PredicateF1=24.2

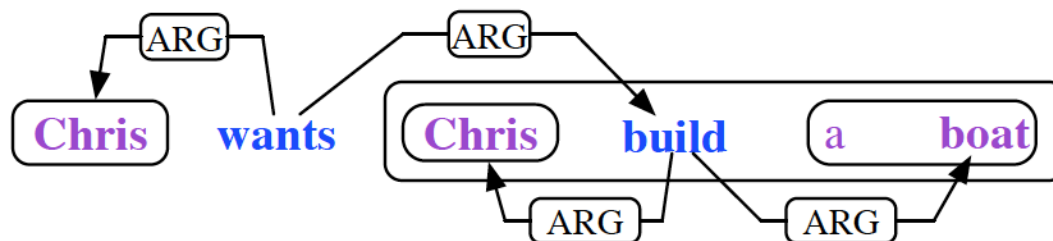
JOINT w/ Moses: BLEU=18.3 / PredicateF1=25.1

JOINT w/ Neural: BLEU=18.9 / PredicateF1=25.8

克里斯想造一艘船。



**Neural
Sequence-to-Sequence
Model**



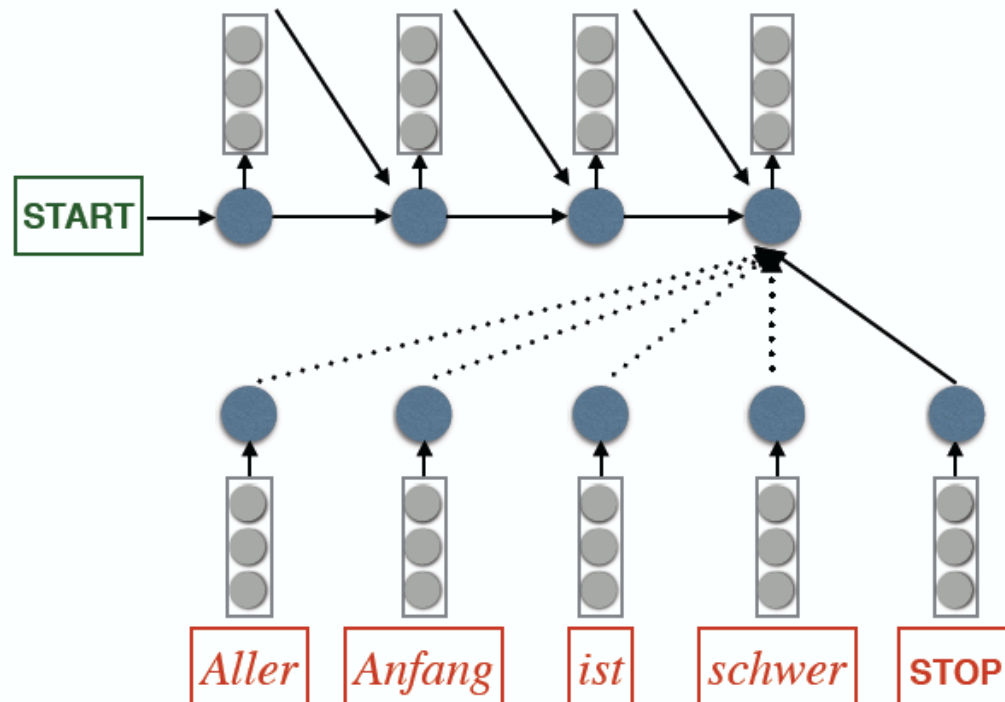
Outline

1. Motivation
2. Problem Definition
3. Pipeline vs. Joint Solution
4. Improved Joint Solution

Sequence generation vs. labeling

- Previously, treat word:label as single token

$[(\text{Chris}:a_h) \text{ wants}:p_h [(\text{Chris}:a_h) \text{ build}:p_h (\text{a}:a \text{ boat}:a_h)]]$



Decompose generation and labeling

$$P(Y, T | X) = \prod_{i=1}^{|Y|} P(y_i, t_i | y_{<i}, t_{<i}, X)$$

target words source words

target labels

previous target words

previous labels

$$= \prod_{i=1}^{|Y|} P(y_i | y_{<i}, t_{\leq i}, X) P(t_i | y_{<i}, t_{<i}, X)$$

Decoder depends on t_i

Predict label t_i

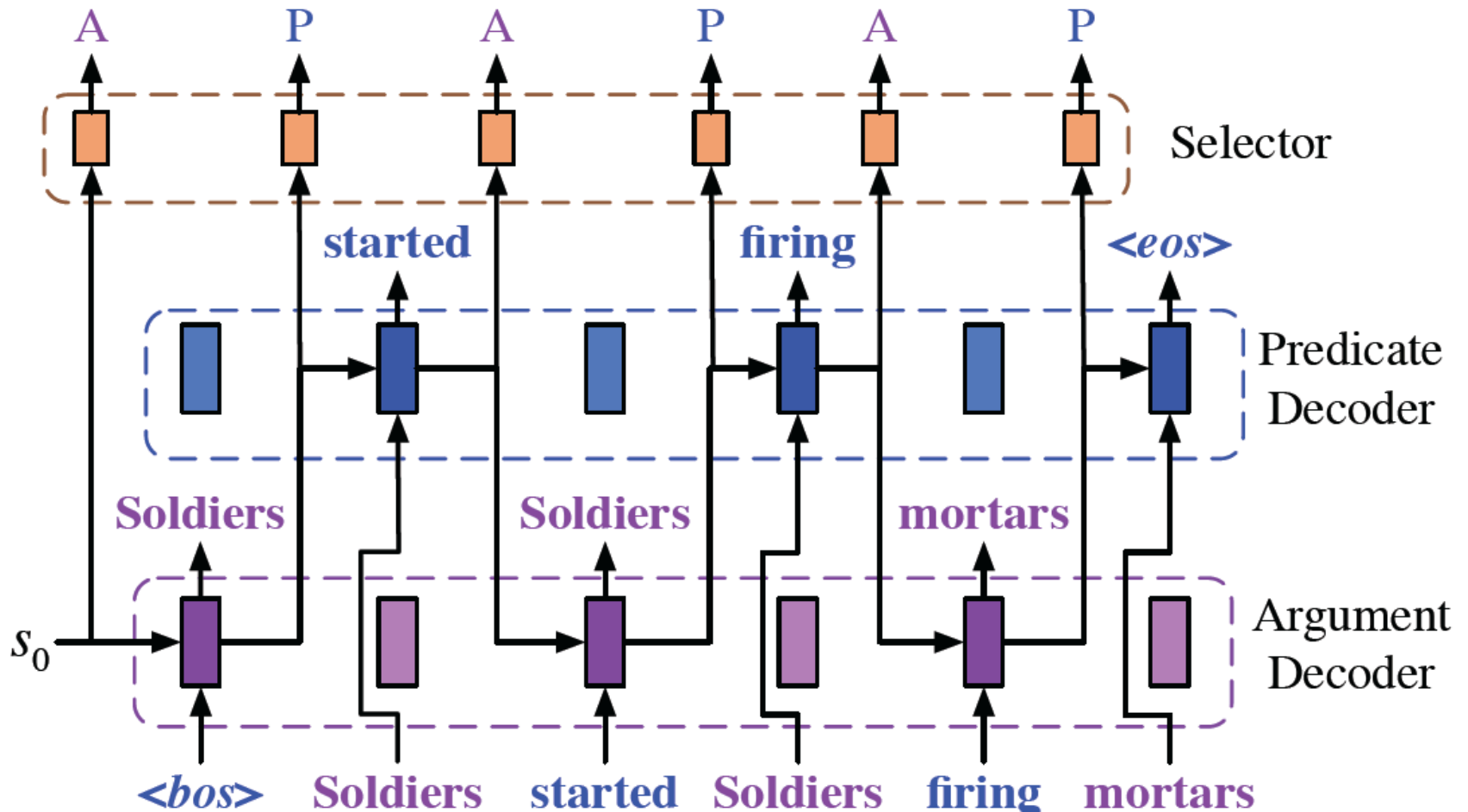
Decompose generation and labeling

$$= \prod_{i=1}^{|Y|} P(y_i \mid y_{<i}, t_{\leq i}, X) P(t_i \mid y_{<i}, t_{<i}, X)$$

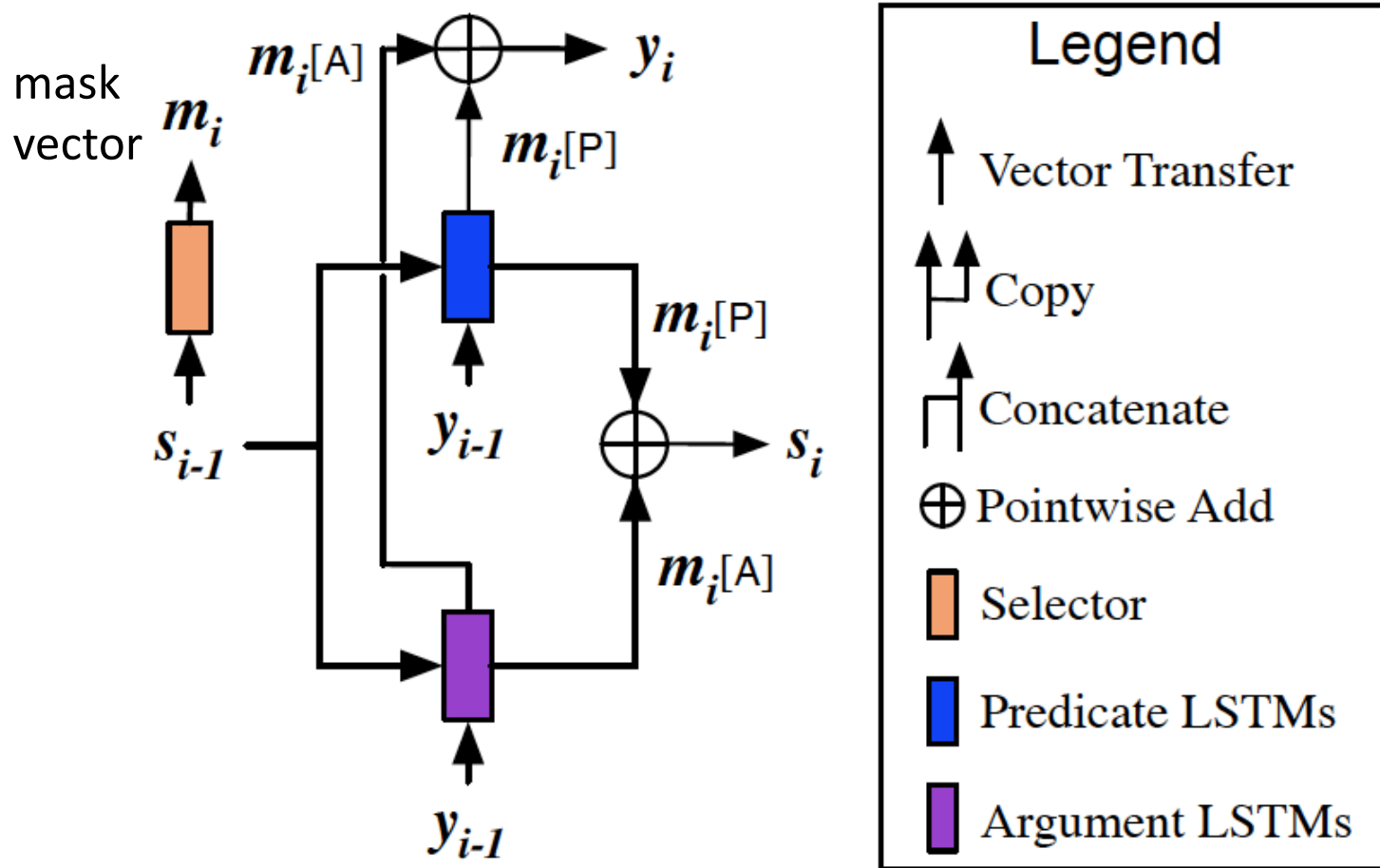
↑ Decoder depends on t_i ↑ Predict label t_i

- Limits increase of target vocabulary
- Models generation process separately by type
 - Given previous word “wanted”:
 - predicate decoder generates “to”, “by”
 - argument decoder generates “a”, “him”

Selective Decoding



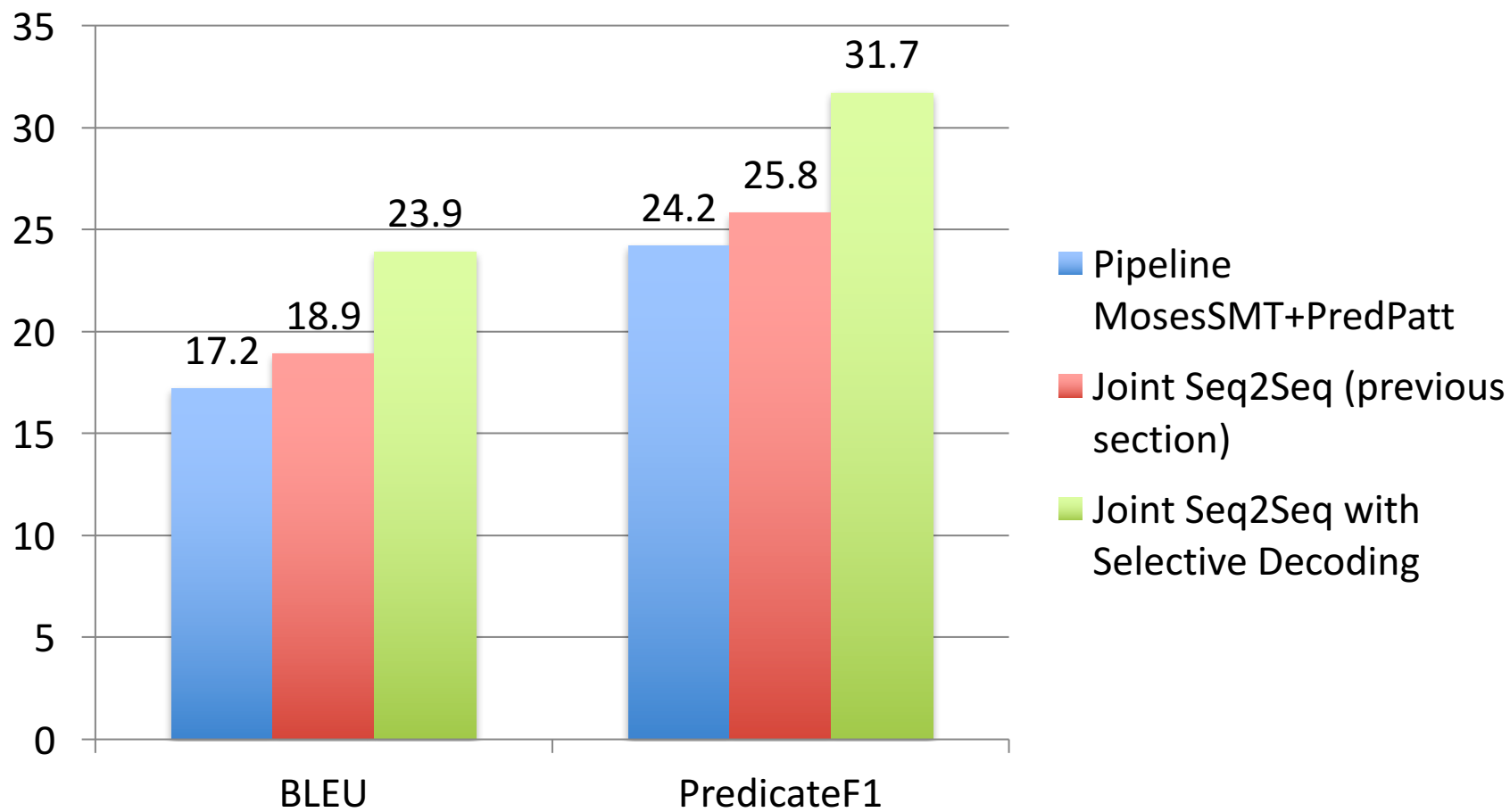
In detail: at each decoding step



Attention context (c_i)
is omitted in the figure

$$s_i = \sum_{t_i \in \mathcal{T}} m_i[t_i] f_{t_i}(y_{i-1}, s_{i-1}, c_i)$$

Results on Chinese-English task



If we only care about MT (not MT+IE)

克里斯想造一艘船。

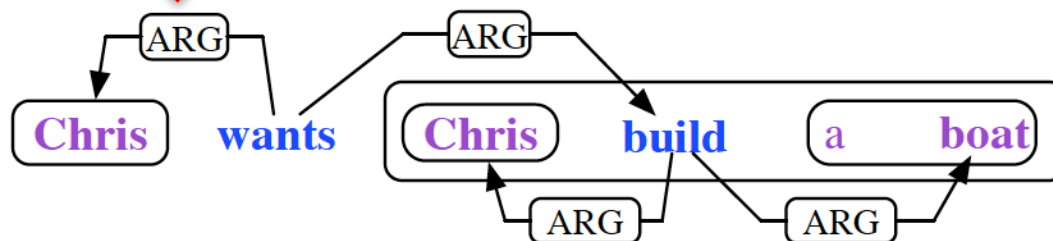
Standard seq2seq

BLEU = 24.92

Chris wants to build a boat

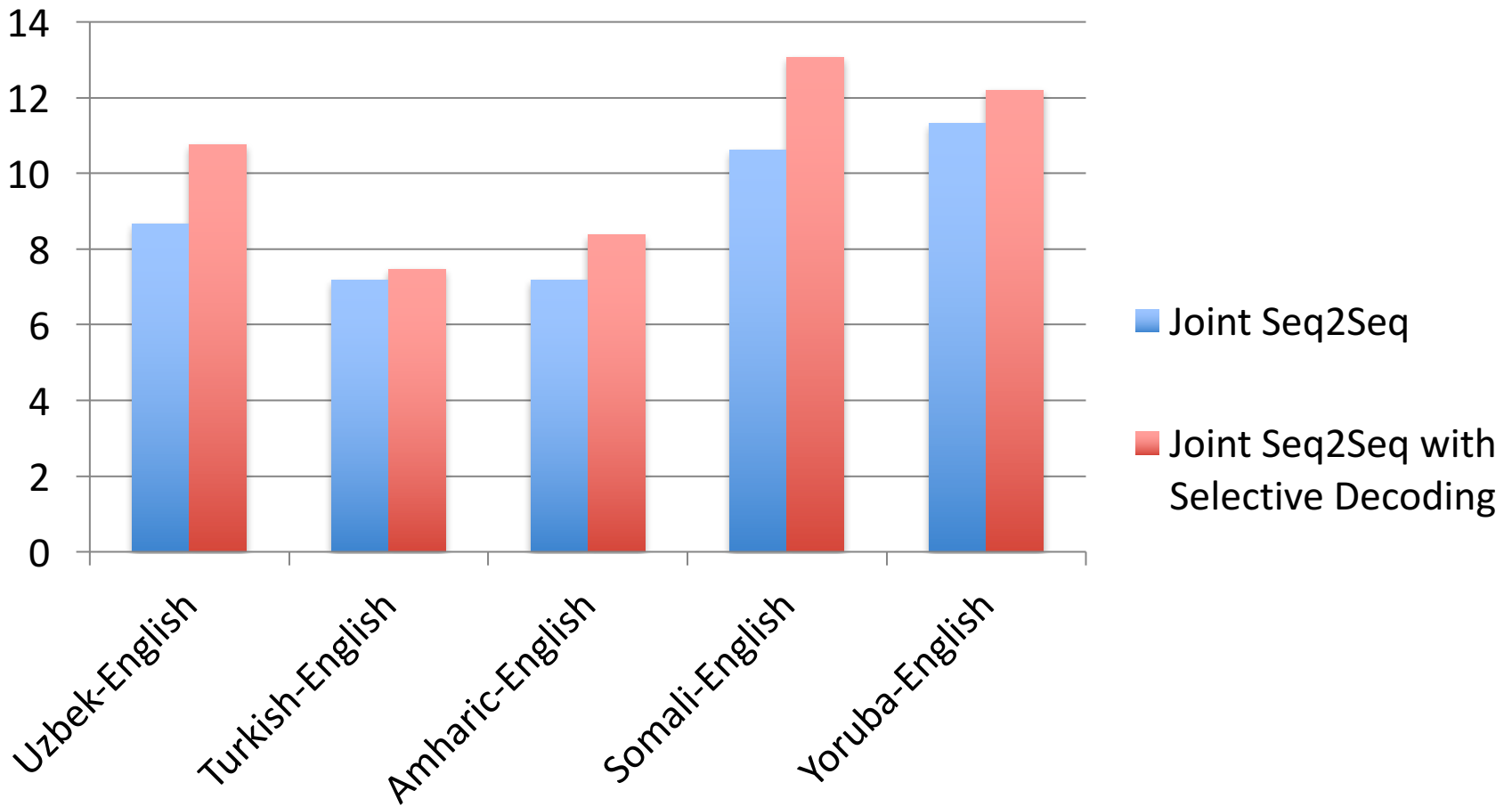
Seq2Seq with selective decoding,
then ignore labels

BLEU = 25.16



BLEU on Low-Resource Languages

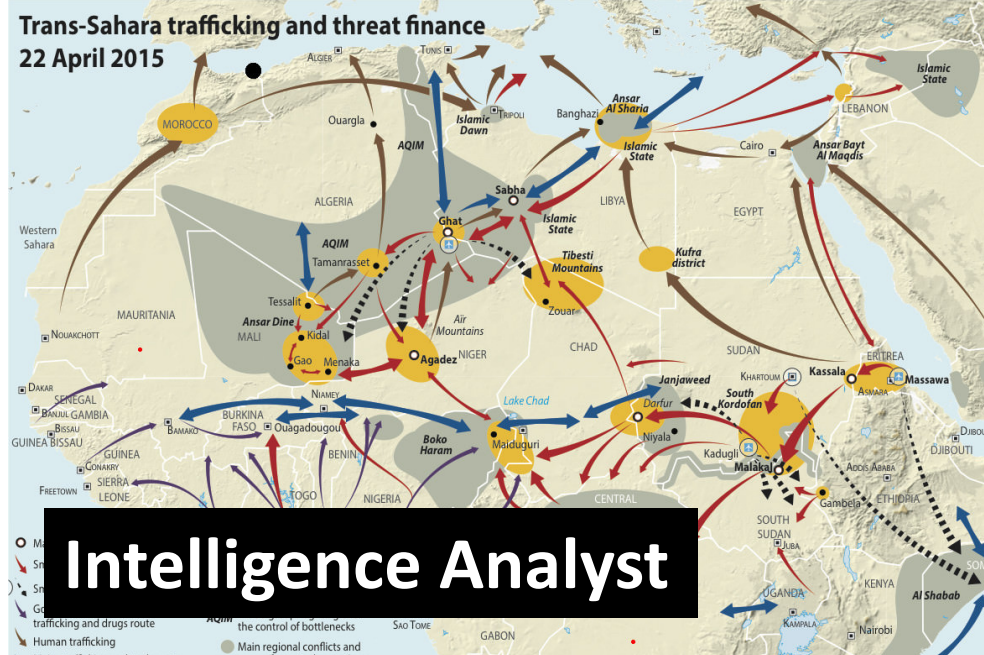
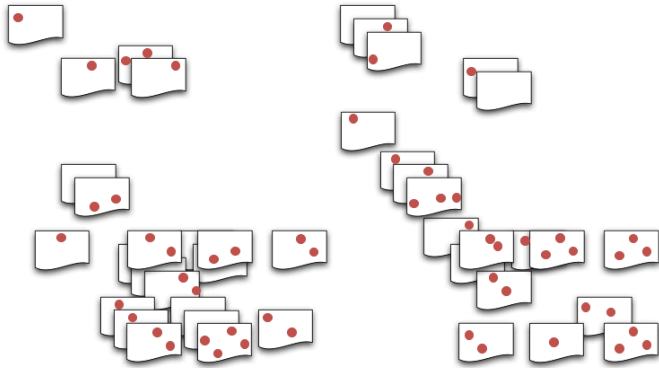
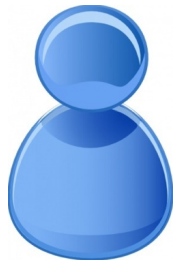
(Data from DARPA LORELEI Project)



NOT SHOWN: The winner is less clear for Predicate F1 (8-14%) 43

Summary

Support users with complex information needs

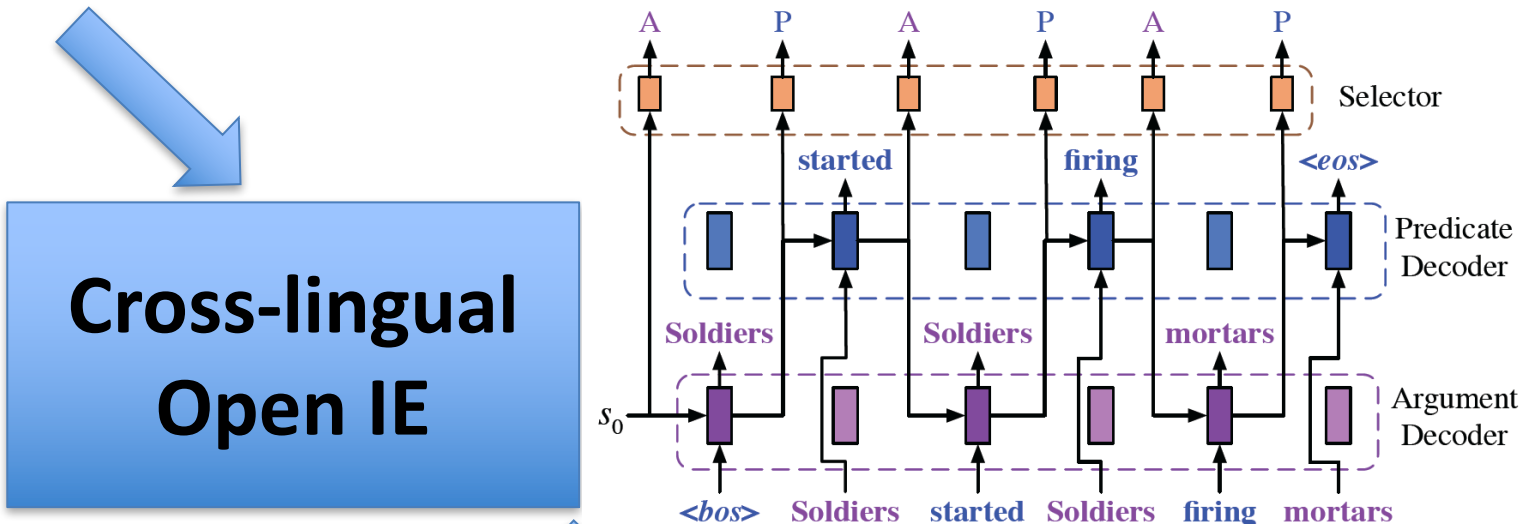


Aid Worker

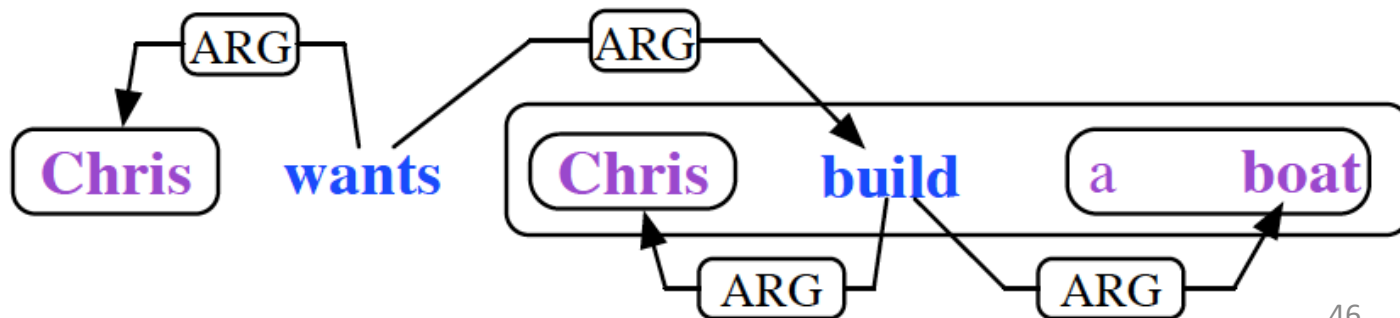


Input: Chinese sentence

克里斯想造一艘船。



Output: English tuples, e.g. Relation(arg1,arg2)



Next Steps

- Integration with analyst search engine
- Directly optimize IE objective, not likelihood
- Explore Selective Decoding for other problems

Thanks!

- To Learn More:
 - S. Zhang, K. Duh, B. Van Durme. “MT/IE: Cross-lingual Open Information Extraction” (EACL2017)
 - S. Zhang, K. Duh, B. Van Durme. “Selective Decoding for Cross-lingual Open Information Extraction” (IJCNLP2017)
 - Code on GitHub