# Semantic Role Labeling for Learner Chinese: the Importance of Syntactic Analysis and L2-L1 Parallel Data

**Zi Lin**, Yuguang Duan, Yuanyuan Zhao, Weiwei Sun, Xiaojun Wan

Peking University

{*zi.lin, ariaduan, zhao_yy, ws, wanxiaojun*}*@pku.edu.cn*
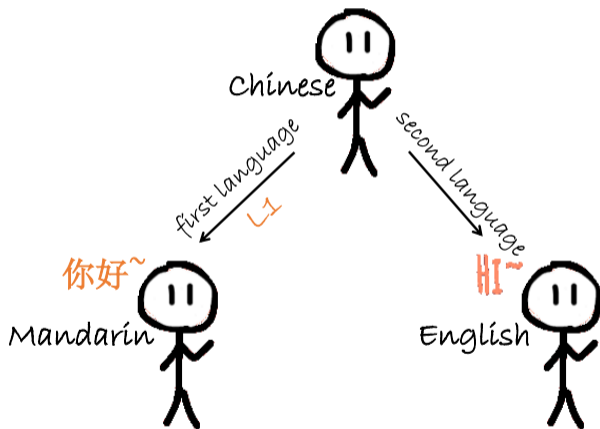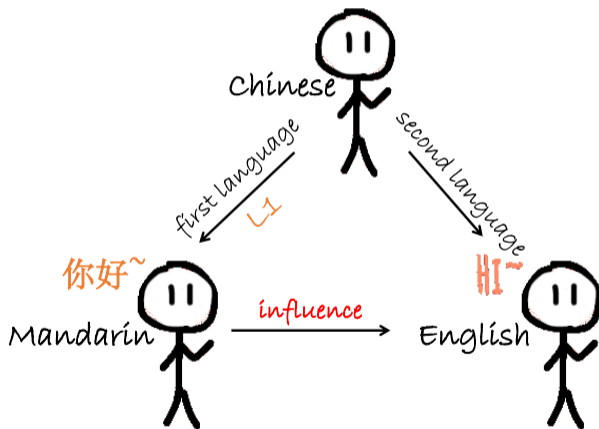
October 25, 2018

# Overview

# Outline

# What is interlanguage?

A second language (or L2) which preserves some features of their first language (or L1).

# What is interlanguage?

A second language (or L2) which preserves some features of their first language (or L1).

# What is interlanguage?

A second language (or L2) which preserves some features of their first language (or L1).

# Interlanguage is everywhere…

Interlanguage is everywhere...



## Social Network

**Semantic Role Labeling for Learner Chinese:**
**the Importance of Syntactic Parsing and L2-L1 Parallel Data**

Zi Lin[123], Yuguang Duan[3], Yuanyuan Zhao[125], Weiwei Sun[124] and Xiaojun Wan[12]

[1]Institute of Computer Science and Technology, Peking University
[2]The MOE Key Laboratory of Computational Linguistics, Peking University
[3]Department of Chinese Language and Literature, Peking University
[4]Center for Chinese Linguistics, Peking University
[5]Academy for Advanced Interdisciplinary Studies, Peking University
{zi.lin,ariaduan,ws,wanxiaojun}@pku.edu.cn, zhaoyy1461@gmail.com

# And perhaps your paper...

PEKING UNIVERSITY

# Outline

北京大学
PEKING UNIVERSITY

# L2-L1 Parallel Data

Collect a large dataset of L2-L1 parallel texts of **Mandarin** by exploring "language exchange" social networking services – lang-8[1].



Post in the language that you are learning.

Native speakers correct your writing!

[1] http://lang-8.com/

# Data for SRL annotation

Initial collection 1,108,907 pairs → clean up → 717,241 pairs → manual selection → 600 pairs → segmentation → SRL annotation

4 typologically different mother tongues

# Data for SRL annotation



Initial collection 1,108,907 pairs → clean up → 717,241 pairs → manual selection → 600 pairs → segmentation → SRL annotation

4 typologically different mother tongues

# Data for SRL annotation



Initial collection 1,108,907 pairs → clean up → 717,241 pairs → manual selection → 600 pairs → segmentation → SRL annotation

4 typologically different mother tongues

# Data for SRL annotation

# Data for SRL annotation



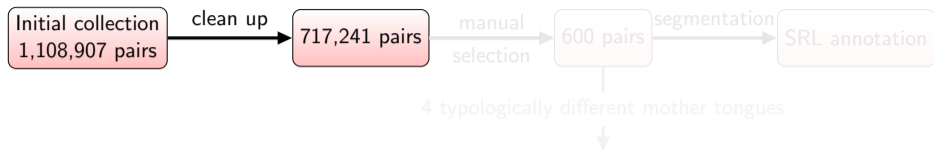| Language | Family |
|----------|--------------|
| Chinese | Sino-Tibetan |
| Russian | Slavic |
| Arabic | Semitic |
| Japanese | Unknown |
| English | Germanic |

# Two Questions

1. Can human understand interlanguage robustly?

2. Can automatic system produce high-quality semantic structures?

# Can human understand interlanguage robustly?

☹ It is difficult to define the syntactic formulism of learner language.

☺ But sometimes we can understand what they mean...

Why not Semantics?

# Can human understand interlanguage robustly?

☹ It is difficult to define the syntactic formulism of learner language.

☺ But sometimes we can understand what they mean...



Why not Semantics?

# Can human understand interlanguage robustly?

☹ It is difficult to define the syntactic formulism of learner language.

☺ But sometimes we can understand what they mean...



Why not Semantics?

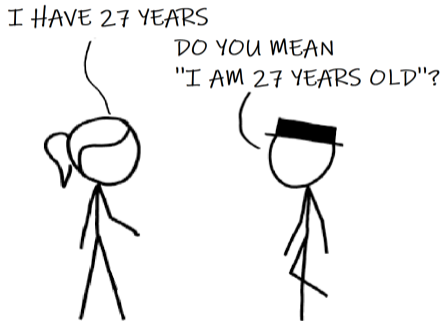# Semantic Role Labeling

**Argument (AN):** | Who | did | what | to | whom |?
**Adjunct (AM):** | When |, | where |, | why | and | how |?

# Inter-annotator agreement

- **Annotator:** two Linguistic students
- **The first 50-sentence trial set:** adapting and refining CPB secification
- **The rest 100-sentence set:** reporting the inter-annotator agreement



Inter-annotator agreement

# Inter-annotator agreement

- **Annotator:** two Linguistic students
- **The first 50-sentence trial set:** adapting and refining CPB secification
- **The rest 100-sentence set:** reporting the inter-annotator agreement



Inter-annotator agreement

# Two Questions

1. Can human understand interlanguage robustly?

2. Can automatic system produce high-quality semantic structures?

# Outline

北京大学
PEKING UNIVERSITY

# Three SRL systems

- The Necessity of Parsing for Predicate Argument Recognition. (2002). Gildea and Palmer.
- Semantic Role Labeling Using Different Syntactic Views (2005). Pradhan et al.
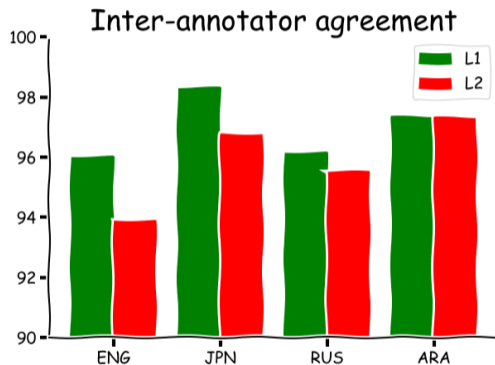- Syntax for Semantic Role Labeling, To Be, Or Not To Be. (2018). He et al.
- Linguistically-Informed Self-Attention for Semantic Role Labeling. (2018). Strubell et al. EMNLP 2018 Best Paper



Trained on Chinese TreeBank that has SRL in CPB

**Parsers**

| Berkeley parser | Performance < | Minimal span-based parser |

**Systems**

| PCFGLA-parser-based SRL system | Neural-parser-based SRL system | Neural syntax-agnostic SRL system |

Trained on Chinese PropBank (CPB)

PEKING UNIVERSITY

# Results



Performance on L1 & L2

A: PCFGLA-parser-based
B: Neural-parser-based
C: Neural syntax-agnostic

# Results



Performance on L1 & L2

A: PCFGLA-parser-based
B: Neural-parser-based
C: Neural syntax-agnostic

# Findings



Performance on L1 & L2

A: PCFGLA-parser-based
B: Neural-parser-based
C: Neural syntax-agnostic

L1
L2
L1-C

The syntax-based systems are more robust when handling learner texts.

# Findings



Performance on L1 & L2

A: PCFGLA-parser-based
B: Neural-parser-based
C: Neural syntax-agnostic

Legend: L1 (green), L2 (red), L2-B (black)

Categories: ENG, JPN, RUS, ARA, ALL

The better the parsing results we get, the better the performance on L2 we achieve.

# Outline

# Why syntactic analysis is important?



用 汉语 也 说话 快　对我来说　很　难 啊。
Using Chinese  also speaking quickly      for  me    very   hard.

| A0 | | AM | AM | rel |

Gold

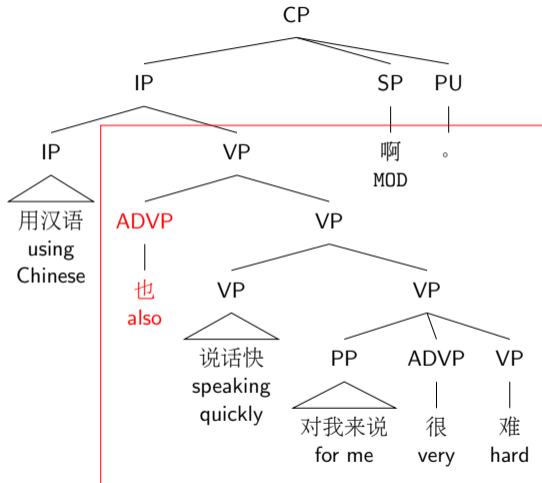| A0 | AM | | AM | AM | rel |

Syntax-based system

| A0 | | AM | rel |

Neural end-to-end system

Using Chinese and also speaking quickly is very hard for me

# Why syntactic analysis is important?
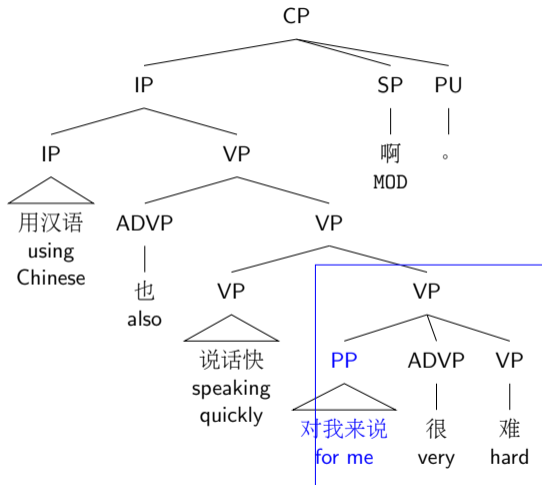


▶ Though the whole structure is ill-formed

# Why syntactic analysis is important?



- Partial of the sentence can be well-formed.

# A new Questions

1. Can human understand interlanguage robustly?

2. Can automatic system produce high-quality semantic structures?

↓

3. Can we improve the SRL performance on interlanguage?

北京大学
PEKING UNIVERSITY

# Outline

北京大学
PEKING UNIVERSITY

# Leveraging L2-L1 Parallel Data

☺ 我　喜欢　做　　中国菜
　 I　 like　cooking　Chinese food

☺ 我　喜欢　做饭
　 I　 like　cooking meal

☹ 我　喜欢　做饭　　中国菜
　 I　 like　cook-meal　Chinese food

# Leveraging L2-L1 Parallel Data

☺ 我　喜欢　做　　　中国菜
　　I　 like 　cooking 　Chinese food

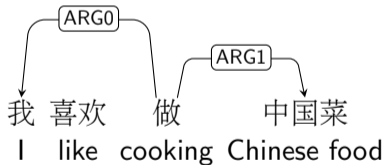☺ 我　喜欢　做饭
　　I　 like　 cooking meal

☹ 我　喜欢　做饭　　中国菜
　　I　 like　 cook-meal　 Chinese food

# Leveraging L2-L1 Parallel Data

:) 我　喜欢　做　　　中国菜
　　I　　like　cooking　Chinese food

:) 我　喜欢　做饭
　　I　　like　cooking meal

:( 我　喜欢　做饭　　中国菜
　　I　　like　cook-meal　Chinese food

# Leveraging L2-L1 Parallel Data

$\langle predicate, argument, role \rangle$ tuples



L1:

ARG0     ARG1

我 喜欢     做中国菜

I   like   cooking Chinese food

ARG0

ARG1

我 喜欢 做 中国菜

I   like   cooking Chinese food

L2:

ARG0     ARG1

我 喜欢     做饭中国菜

I   like   cook-meal Chinese food

ARG0

ARG1

我 喜欢 做饭 中国菜

I   like   cook-meal Chinese food

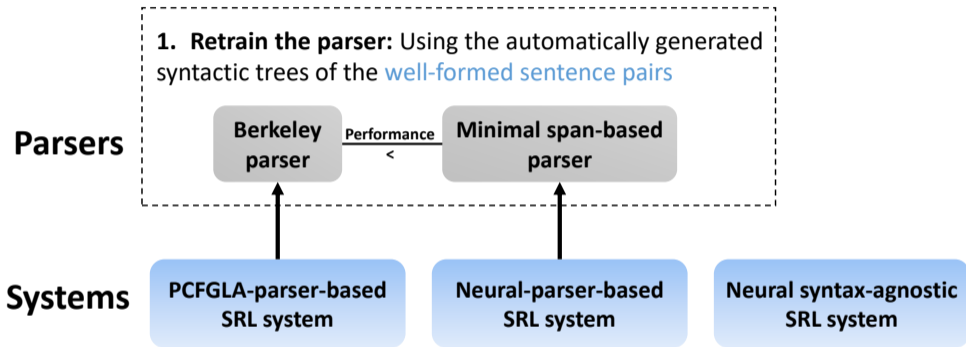\# of shared tuples = 1

# Leveraging L2-L1 Parallel Data

Metric for comparing SRL results

- L2-recall:
  (# of shared tuples) / (# of tuples of the result in L2)
- L1-recall:
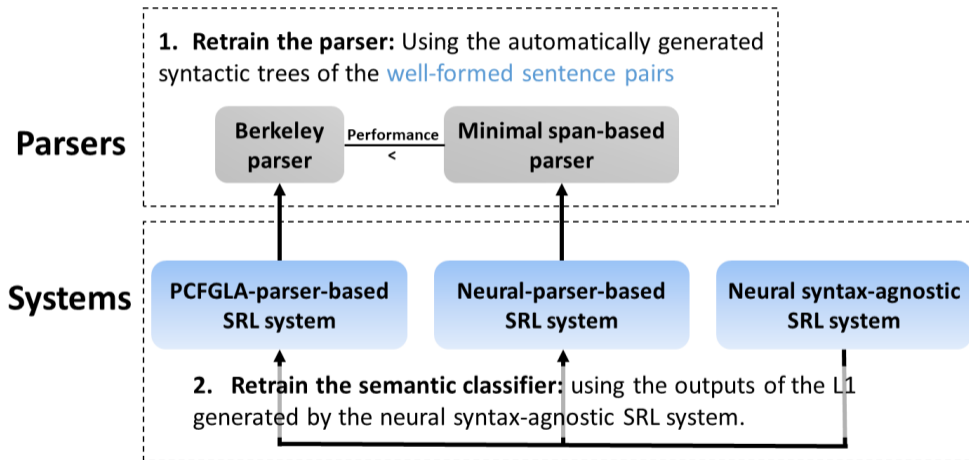  (# of shared tuples) / (# of tuples of the result in L1)
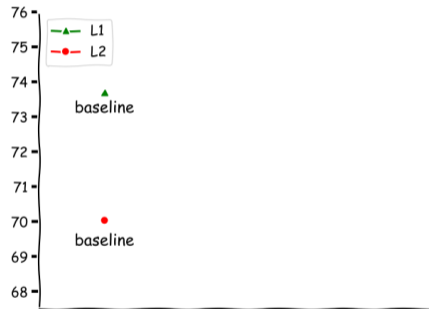
Well-formed sentence pair if both are greater than $\lambda$

# Retraining two essential modules

**Parsers**

**1. Retrain the parser:** Using the automatically generated syntactic trees of the well-formed sentence pairs

| Berkeley parser | Performance < | Minimal span-based parser |
|---|---|---|

**Systems**

| PCFGLA-parser-based SRL system | Neural-parser-based SRL system | Neural syntax-agnostic SRL system |
|---|---|---|

北京大学
PEKING UNIVERSITY

# Retraining two essential modules



**Parsers**

1. **Retrain the parser:** Using the automatically generated syntactic trees of the well-formed sentence pairs

Berkeley parser — Performance < — Minimal span-based parser

**Systems**

PCFGLA-parser-based SRL system

Neural-parser-based SRL system

Neural syntax-agnostic SRL system

2. **Retrain the semantic classifier:** using the outputs of the L1 generated by the neural syntax-agnostic SRL system.

北京大学
PEKING UNIVERSITY
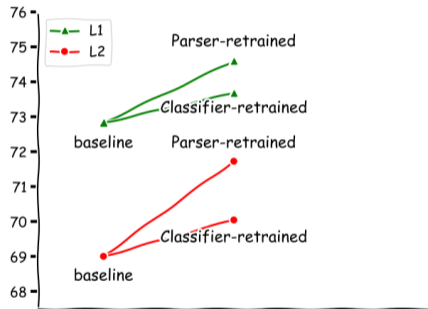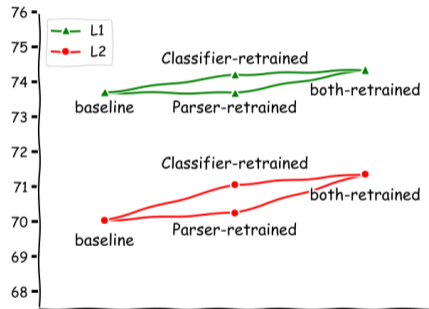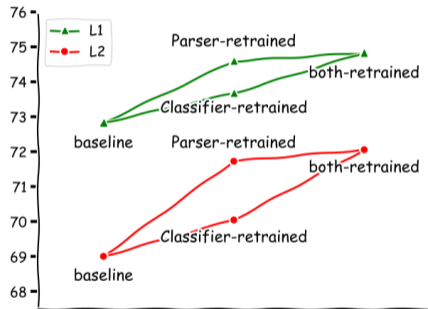
# Results

# Results

# Results

# Thanks for your attention!

**Zi Lin** is planning to apply for PhD program in CS or linguistics this fall. Email me at `zi.lin@pku.edu.cn` if you are interested!

北京大学
PEKING UNIVERSITY