

# Learning Language Structures through Grounding

Freda Shi

Toyota Technological Institute at Chicago

*freda@ttic.edu*

September 21, 2023

@Peking University

## Conversations between Freda and ChatGPT

Here is a sentence: **The cat near the children meows.** How many nouns, noun phrases and verbs are there in the given sentence? What are they?



There are three nouns in the sentence: "cat," "children," and "meows."  
There is one noun phrase: "the cat near the children."  
There is one verb: "meows."

Large language models do not understand syntactic structures well.

Large language models do not understand discourse structures, and can be easily distracted by irrelevant context [Shi et al., ICML 2023].

Can you show me a picture of a cat?



I'm sorry, as an AI language model, I don't have the ability to display images...

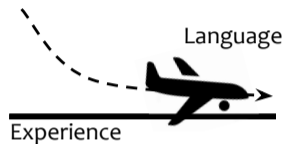
Large language models need

- A better understanding of language structures;
- A better connection with the real world.

Learning **Language Structures** through **Grounding**

# What is Grounding?

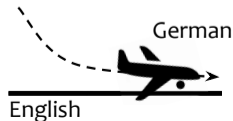
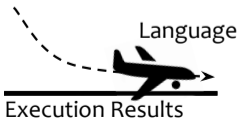
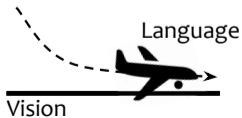
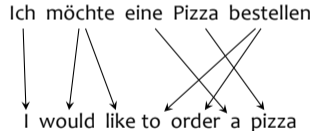
Grounding: Landing the *plane of language* onto the *ground of experience*.



A cat is standing on the lawn.



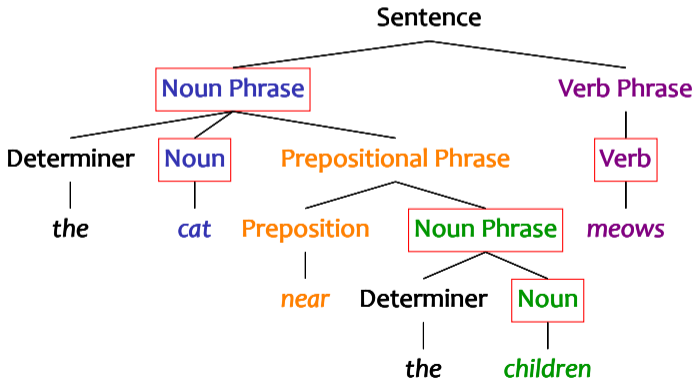
What is the factorial of 5?  
 ↓ Semantic Parser  
`from utils import factorial`  
`print(factorial(5))`  
 ↓ Python Interpreter  
 120



## What are Language Structures?

The **cat** **near** **the** **children** **meows**.

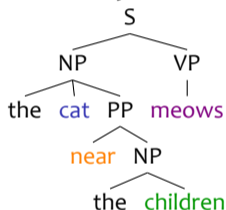
How many nouns, noun phrases and verbs are there in the sentence?  
What are they?



# What are Language Structures?

## The **cat** **near** the **children** **meows**.

### Constituency Parse Tree



### Dependency Parse Tree



### Syntactic Structures

### Truth-Conditional Semantics

$\lambda x. \lambda y. \text{cat}(x) \wedge \text{children}(y) \wedge$   
 $\text{near}(x, y) \wedge \text{meow}(x)$

### SQL

```
SELECT * FROM catsNearChildren
WHERE meows = true;
```

### Python

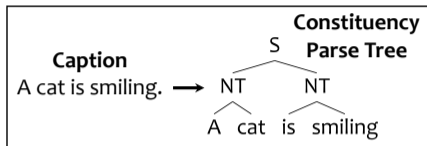
```
def find(cats):
    for cat in cats:
        if cat.near(children) and cat.meows:
            return cat
```

### Semantic Structures

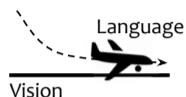
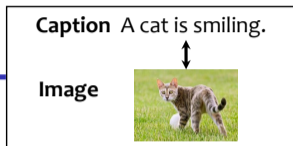
## Learning Language Structures through Grounding

- Why do we care about learning language structures?  
Language structures can
  - Model human language processing;
  - Test or even inform linguistic theories;
  - Enable better interaction between humans and machines;
  - Improve machine learning models.
- Language structures are useful, but expensive to annotate.  
Many grounding signals exist naturally.
- Byproduct: Derived methods and analysis can benefit broader NLP community.

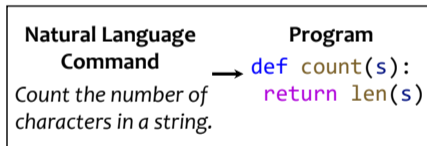
• Learning to parse sentences through visual grounding



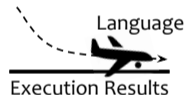
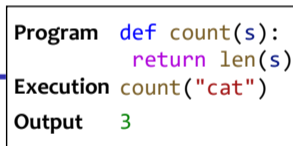
As Supervision



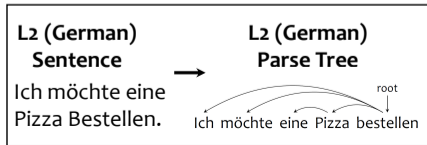
• Learning semantic parses through execution results



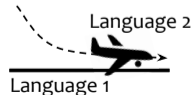
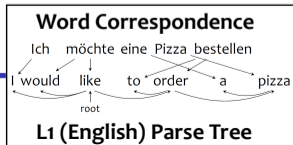
Informed Sampling



• Learning to parse another language through cross-lingual grounding

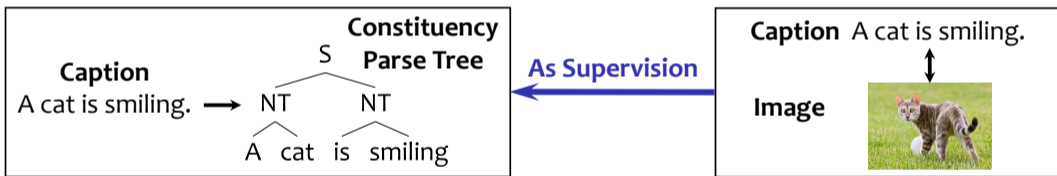


As Supervision





# Part I: Learning to Parse Sentences through Visual Grounding



**Question:** Can visual grounding help induce linguistic structures?

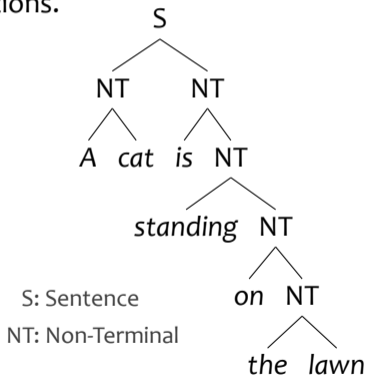
## Problem Formulation

**Task:** Visually grounded grammar induction.

**Input:** Captioned images.

**Output:** Linguistically plausible structure for captions.

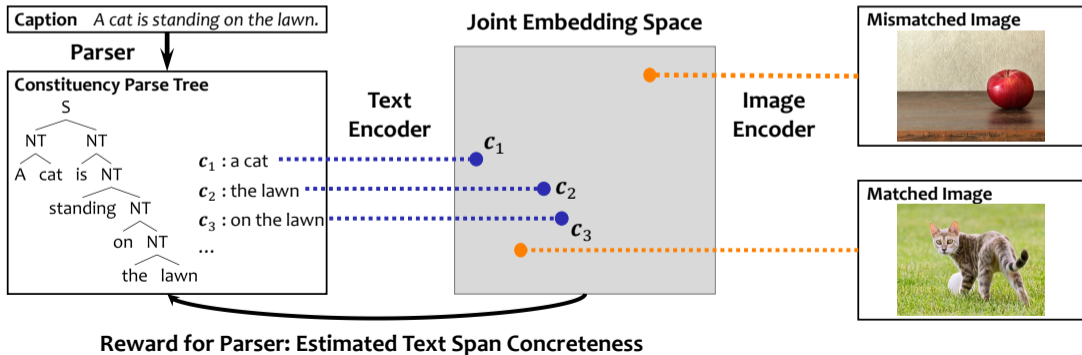
*A cat is standing on the lawn.*



# The Visually Grounded Neural Syntax Learner (VG-NSL)

**Hypothesis:** More visually concrete word spans are more likely to be constituents.

Joint Embedding Space: Higher similarity for matched image-constituent pairs;  
Lower similarity for mismatched pairs.



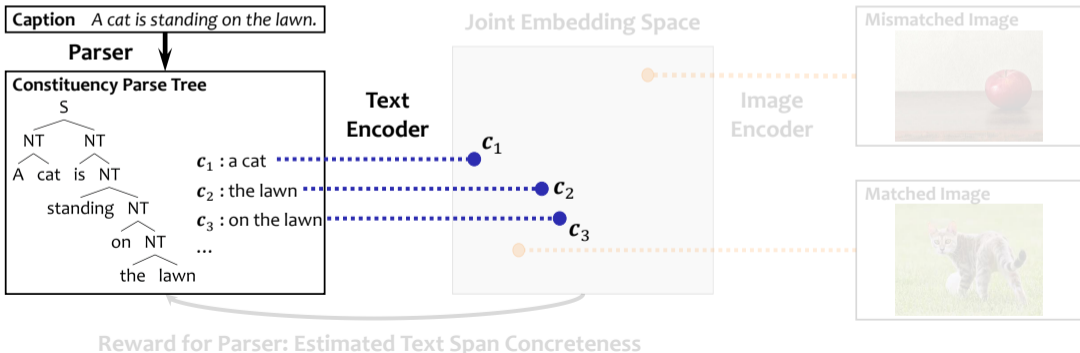
## VG-NSL: Text Parser and Encoder

**Parser Output:** Predicted phrases and their vector representations.

$P_{\Theta}(\mathbf{v}^{\text{"a cat"}})$ : Probability of "a cat" to be a constituent.

$\mathbf{V}$ : Semantic representation of words and word spans.

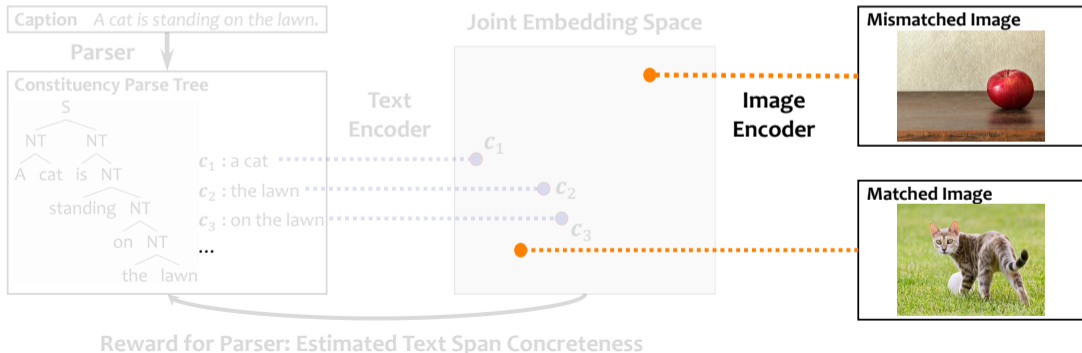
$\Theta$ : Structure of parse trees.



## VG-NSL: Image Encoder

**Image Encoder:** Frozen ResNet (He et al., 2015) + Linear Projection.

$$\mathbf{u}_{img} = \Phi \cdot \text{ResNet}(img)$$

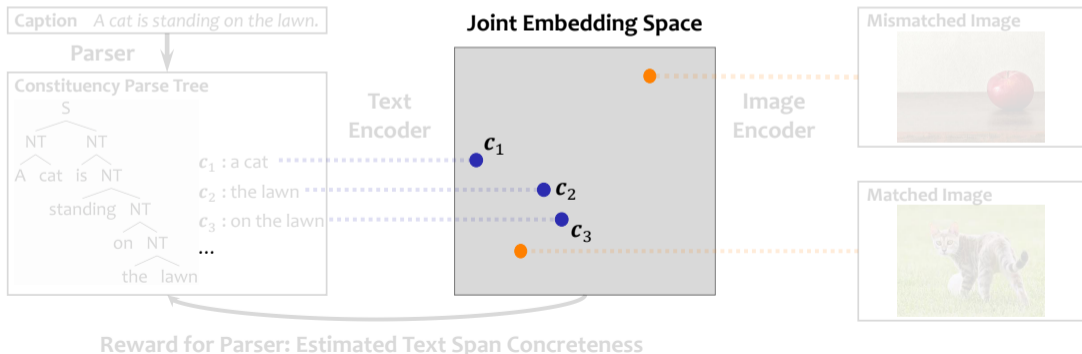


## VG-NSL: Joint Visual-Semantic Embedding Space

**Parameters for Text Encoder:** Word representations  $\mathbf{V}$ , parser parameters  $\Theta$ .

**Parameters for Image Encoder:** Linear projector  $\Phi$ .

**Joint Embedding Space:** Train  $\mathbf{V}$  and  $\Phi$  – align meanings of word spans and images.



## VG-NSL: Joint Embedding Space

**Key Idea:** Higher similarity for matched image-constituent pairs,  
Lower similarity for mismatched pairs.

**Approach:** Minimize hinge-based triplet loss (Kiros et al., 2015)  
between images and constituents.

$$\mathcal{L}(i, c; \mathbf{V}, \Phi) = \sum_{(i', c') \neq (i, c)} [\text{sim}(i', c) - \text{sim}(i, c) + \delta]_+ + [\text{sim}(i, c') - \text{sim}(i, c) + \delta]_+$$



$\text{sim}(\cdot, \cdot) = \cos(\cdot, \cdot)$      $[\cdot]_+ = \max(0, \cdot)$      $\delta$  : margin score

## VG-NSL: Quantify Visual Concreteness

Joint Embedding Space: Higher similarity for matched image-constituent pairs;  
Lower similarity for mismatched pairs.

Image  $i$ Candidate  
Constituent  $c$ 

$a\ cat$   
 $on\ the$

Another Image  $i'$ 

$$\ell(c; i, i') = \text{sim}(i', c) - \text{sim}(i, c)$$

**Value of  $\ell$**

$$\text{sim}\left(\text{img}(\text{apple}), a\ cat\right) = 0.2 \quad \text{sim}\left(\text{img}(\text{cat}), a\ cat\right) = 0.9$$

$$\ell = -0.7$$

$$\text{sim}\left(\text{img}(\text{apple}), on\ the\right) = 0.4 \quad \text{sim}\left(\text{img}(\text{cat}), on\ the\right) = 0.4$$

$$\ell = 0$$

**Key Idea:** Smaller  $\ell(c)$   $\iff$   $c$  is more visually concrete.

Quantify *visual concreteness* of word spans using loss values.

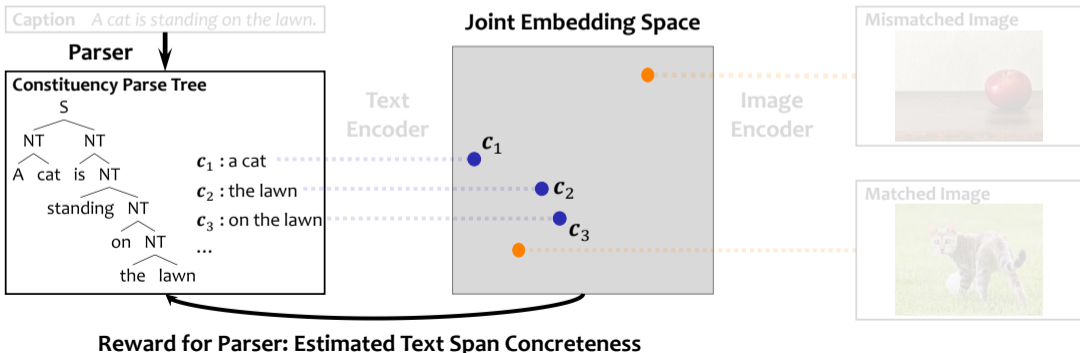


## VG-NSL: Concreteness as Rewards for Text Parser

REINFORCE (Williams, 1992) as the gradient estimator for parsing parameter  $\Theta$ :

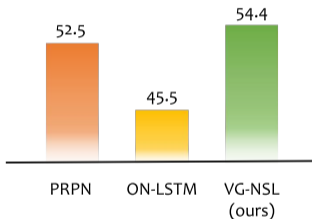
$$\Theta \leftarrow \Theta + \eta \cdot \nabla_{\Theta} \sum_{(i,c)} p_{\Theta}(c) \text{concreteness}(c; i)$$

$\eta$  : learning rate



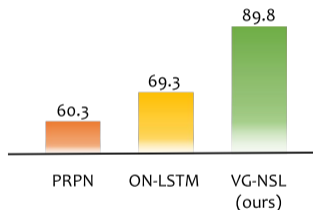
## VG-NSL: Results on the MSCOCO (Lin et al., 2014) Dataset

### F<sub>1</sub> Score (↑)



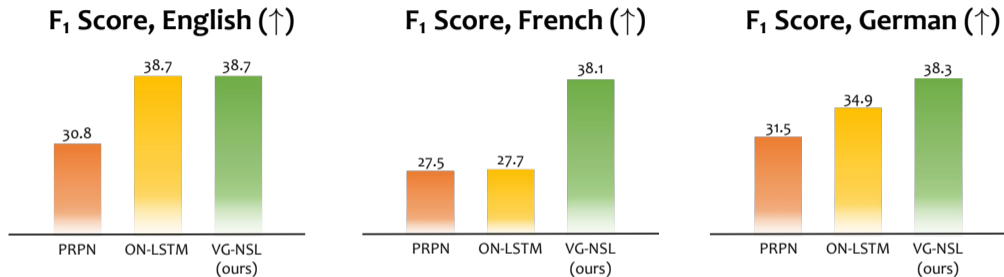
Text-Only Models : { PRPN: Shen et al., 2018  
ON-LSTM: Shen et al., 2019

### Self-F<sub>1</sub> Score (↑)



5 Runs: { Same hyperparameters  
Different random seeds

## VG-NSL: Results on the Multiz0K (Elliott et al., 2016) Dataset



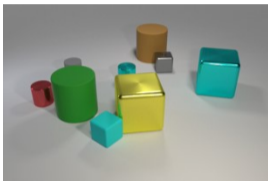
Text-Only Models : { PRPN: Shen et al., 2018  
ON-LSTM: Shen et al., 2019

**Question:** Can visual grounding help induce linguistic structures?

**Answer:** Yes, on syntactic (constituency) parsing.

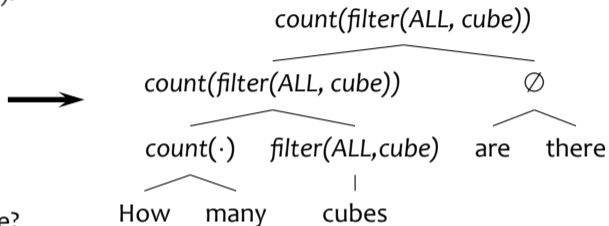
# Joint Syntax and Semantics Induction through Visual Grounding

Dataset: CLEVR (Johnson et al., 2017).



Question: How many cubes are there?

Answer: 4.



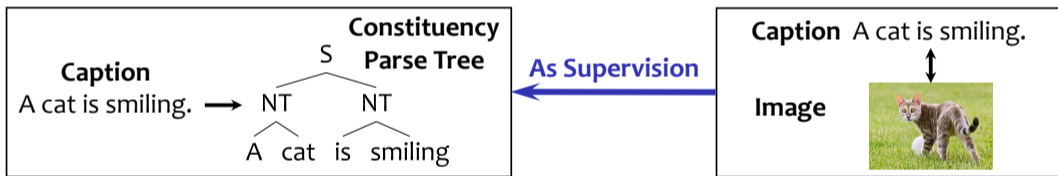
Question answering accuracy ( $\uparrow$ ) on program-depth generalization:

81.6 (prior SotA)  $\rightarrow$  **98.5**

**Question:** Can visual grounding help induce linguistic structures?

**Answer:** Yes, on semantic parsing.

# Part I: Learning Syntactic/Semantic Parses through Visual Grounding



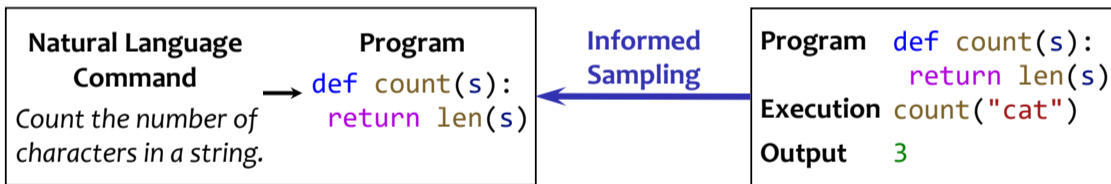
**Question:** Can visual grounding help induce linguistic structures?

**Approach:** Propose the task of visually grounded grammar induction.

**Answer:** Yes, on syntactic (constituency) parsing [SMGL, ACL 2019].

**Answer:** Yes, on semantic parsing [MSWLT, NeurIPS 2021].

## Part II: Learning Semantic Structures through Execution



**Question:** Can execution results, as grounding signals, help learn semantic structures?

**Answer:** Yes, on semantic parsing without program supervision [MSWLT, NeurIPS 2021].

## Problem Formulation

**Task:** Convert natural language to code, leveraging execution (grounding) of programs.

**Input:** Command in natural language.

**Output:** Corresponding program.

**Example Input:** *Write a Python function that counts lowercase letters in a string.*

**Example Output:**

```
def count(string):
    cnt = 0
    for ch in string:
        if ch.islower():
            cnt += 1
    return cnt

import collections
def count(s):
    cnt = collections.Counter(s)
    return sum(
        cnt[c] for c in cnt
        if c.islower()
    )

def count(s):
    return len([
        c for c in s
        if c.islower()
    ])
```

## Background: Codex

Transformer-based generative model for code (Chen et al., 2021).

- **Training:** Model probability of natural language and GitHub code snippets.

$$\max_{\Theta} \prod_{\mathbf{x}} P_{\Theta}(\mathbf{x})$$

$\Theta$  : Model parameters.

$\mathbf{x}$  : Training example.

$$P_{\Theta}(\mathbf{x}) = \prod_{i=1}^{|\mathbf{x}|} P_{\Theta}(x_i | x_1, \dots, x_{i-1})$$

$x_i$  :  $i^{\text{th}}$  token of  $\mathbf{x}$ .

- **Inference:** Generate code conditioned on natural language description.

$$P_{\Theta}(\mathbf{x} | x_1, \dots, x_c) = \prod_{i=c+1}^{c+L} P_{\Theta}(x_i | x_1, \dots, x_{i-1})$$

$x_1, \dots, x_c$  : Natural language description.

$L$  : Maximum decoding step.

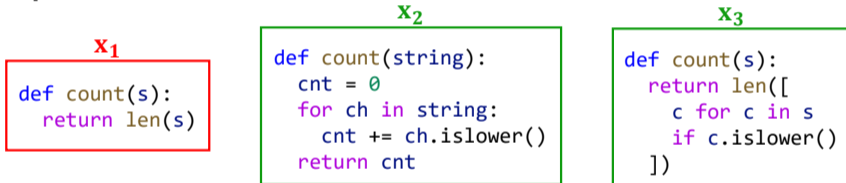


## Natural Language to Code: Decoding Method

**Task:** Generate code conditioned on natural language description.

**Example Input:** Write a Python function that counts lowercase letters in a string.

**Example Output:**



$$P_{\Theta}(x_1 | \dots) = 0.4$$

$$P_{\Theta}(x_2 | \dots) = 0.3$$

$$P_{\Theta}(x_3 | \dots) = 0.3$$

$s_i$ : Execution results of  $x_i$ .



$$P_{\Theta}(s_1 | \dots) = 0.4$$

$$P_{\Theta}(s_2 | \dots) = P_{\Theta}(s_3 | \dots) = 0.3 + 0.3 = 0.6$$

**Key Idea:** Consider program semantics (i.e., execution results)–based equivalent classes.

## Empirical Solution

**Hypothesis:** Codex assigns higher probability to correct execution results.

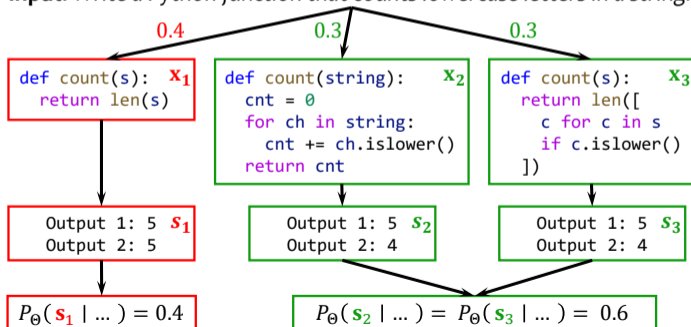
**Approach:** Rank programs with  $P_{\Theta}(\mathbf{s} \mid \dots)$ .

$$P_{\Theta}(\mathbf{s} \mid \dots) = \sum_{\mathbf{x}} P_{\Theta}(\mathbf{x}, \mathbf{s} \mid \dots) = \sum_{\text{exec}(\mathbf{x})=\mathbf{s}} P_{\Theta}(\mathbf{x} \mid \dots)$$

$\mathbf{x}$ : Program.

$\mathbf{s}$ : Execution results.

**Input:** Write a Python function that counts lowercase letters in a string.



**Step 1:** Sample  $\sim P_{\Theta}(\mathbf{x} \mid \dots)$

**Step 2:** Synthesize input cases

Input 1: "hello"

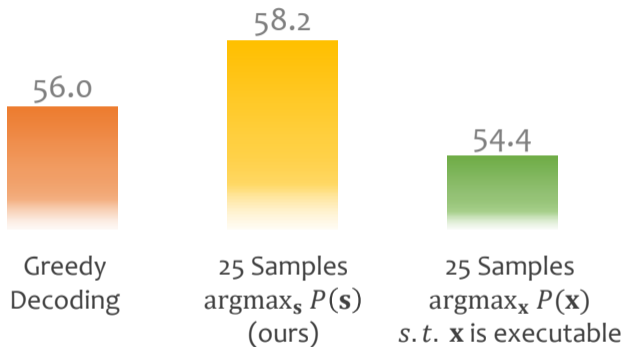
Input 2: "Hello"

**Step 3:** Execute

**Step 4:** Estimate  $P_{\Theta}(\mathbf{s} \mid \dots)$

**Step 5:** Select a program  $\mathbf{x}^*$  with semantics  $\mathbf{s}^* = \arg \max_{\mathbf{s}} P_{\Theta}(\mathbf{s})$  for test

## Results: Natural Language to Python Translation (Austin et al., 2021)



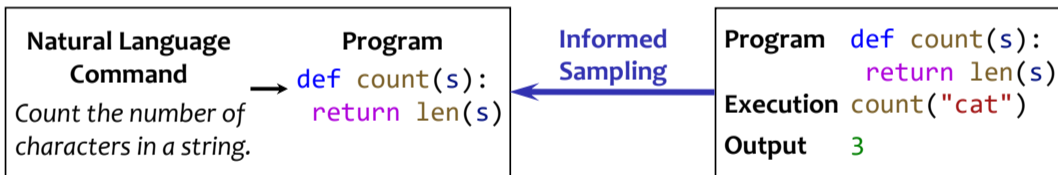
$P(\mathbf{s})$ : Probability of execution results.

$P(\mathbf{x})$ : Probability of program.

**Question:** Can execution results, as grounding signals, help learn semantic structures?

**Answer:** Yes, execution result–based method improves natural language to Python translation.

## Part II: Learning Semantic Structures through Execution

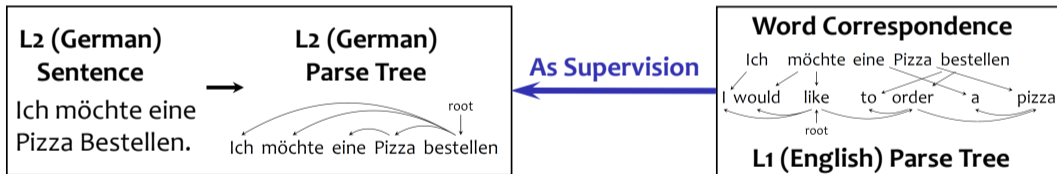


**Question:** Can execution results, as grounding signals, help learn semantic structures?

**Answer:** Yes, on semantic parsing without program supervision [MSWLT, NeurIPS 2021].

**Answer:** Yes, on natural language to code translation [SFGZW, EMNLP 2022].

# Part III: Towards Language-Universal NLP through Cross-Lingual Grounding



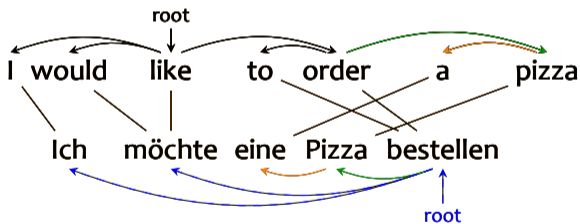
**Question:** Can we transfer NLP models to another language through cross-lingual grounding?

## Problem Formulation

**Task:** Zero-shot cross-lingual dependency parsing.

**Input:** Sentences and dependency parse trees in source language;  
Translated sentences in target language;  
Word correspondence between parallel sentences.

**Output:** Dependency parse trees in target language.



-->: Missing arcs with simple projection

**Key Idea:** Leverage the nature of trained source dependency parser that it can capture “unannotated” dependency relations.

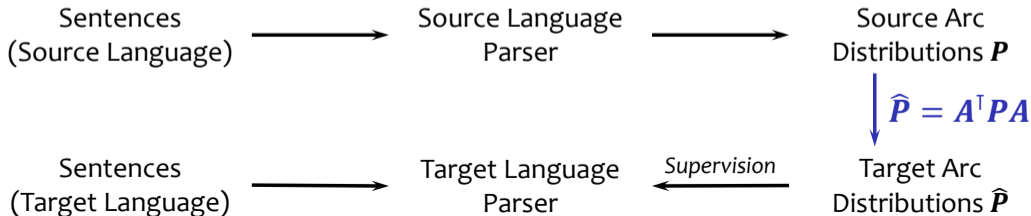
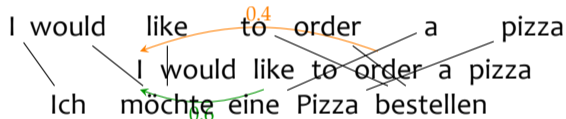
# Target Language Parser Training after Arc Distribution Projection

**Input:** Sentences and dependency parse trees in source language;  
 Translated sentences in target language;  
 Word correspondence between parallel sentences. } Pretrained Multilingual Models

**Output:** Dependency parse trees in target language.

$P_{i,j}$ :  $P(\text{head} = \mathbf{w}_j \mid \text{word} = \mathbf{w}_i)$

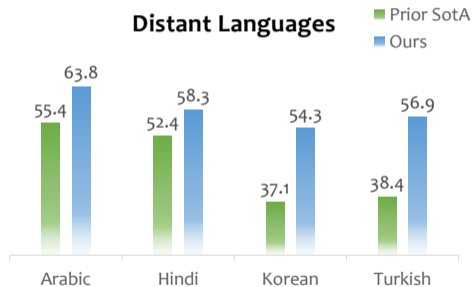
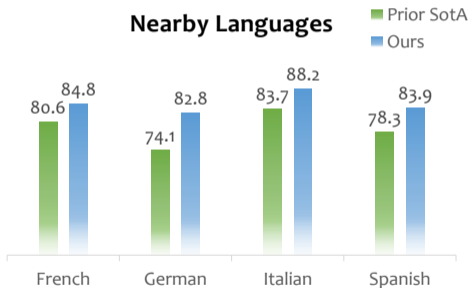
$A_{i,j}$ : 1 if source word  $\mathbf{w}_i$  and target word  $\mathbf{u}_j$  are aligned, otherwise 0.



## Zero-Shot Cross-Lingual Dependency Parsing: Results

**Metric:** Unlabeled attachment score (UAS, ↑).

**Source language:** English (95.8 UAS on English).

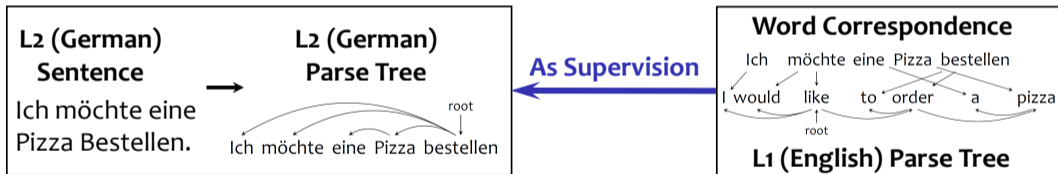


**Question:** Can we transfer NLP models to another language through cross-lingual grounding?

**Answer:** Yes, through substructure (arc) distribution projection.



# Part III: Towards Language-Universal NLP through Cross-Lingual Grounding



**Question:** Can we transfer NLP models to another language through cross-lingual grounding?

- Approach:**
- (1) Train a source language parser;
  - (2) Project the source parser prediction to the target language;
  - (3) Train a target language parser to fit projected distribution.

**Answer:** Yes, through substructure (arc) distribution projection [SGL, ACL 2022].

# Thanks!



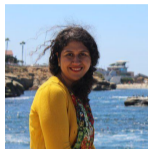
Kevin Gimpel



Karen Livescu



Daniel Fried



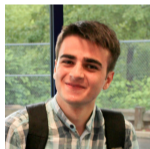
Marjan Ghazvininejad



Roger Levy



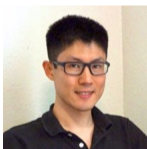
Jiayuan Mao



Mirac Suzgun



Josh Tenenbaum



Sida Wang



Jiajun Wu



Luke Zettlemoyer



Denny Zhou

- Armen Aghajanyan
- Xinyun Chen
- Hyung Won Chung
- David Dohan
- Dipanjan Das

- Markus Freitag
- Lingyu Gao
- Vikram Gupta
- Yuning Jiang
- Mike Lewis

- Lei Li
- Jessy Lin
- Kanishka Misra
- Sebastian Ruder
- Mrinmaya Sachan

- Nathan Scales
- Nathanael Schärli
- Bowen Shi
- Suraj Srivats
- Jian Sun

- Yi Tay
- Shubham Toshniwal
- Soroush Vosoughi
- Eric Wallace
- Xuezhi Wang

- Jason Wei
- Tete Xiao
- Scott Yih
- Ruiqi Zhong
- Hao Zhou