

# ScatterShot: Interactive In-context Example Curation for Text Transformation

Sherry Tongshuang Wu,

*Carnegie Mellon University*

**Hua Shen,**

*PennState University*

Daniel S. Weld,

*University of Washington*

Jeffrey Heer,

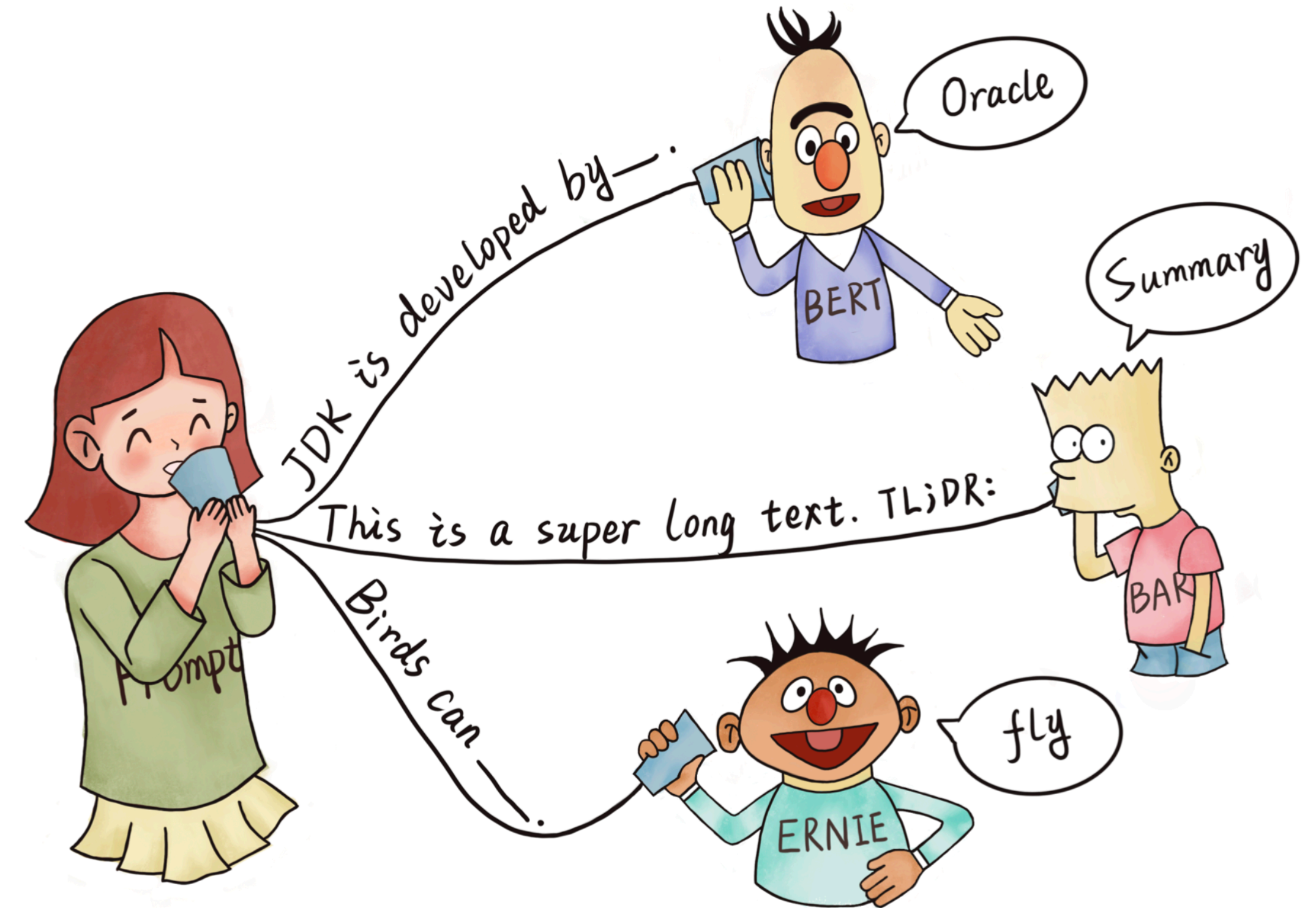
*University of Washington*

Marco Tulio Ribeiro,

*Microsoft*

# What is prompt-based learning with LLMs?

Encourages a **pre-trained** Large Language Model (LLM) to make **particular predictions** by providing a **"prompt"** specifying the task to be done.



# What is prompt-based learning with LLMs?

Encourages a **pre-trained** Large Language Model (LLM) to make **particular predictions** by providing a **"prompt"** specifying the task to be done.

## Prompt Design

In-context Learning

Prompt Search

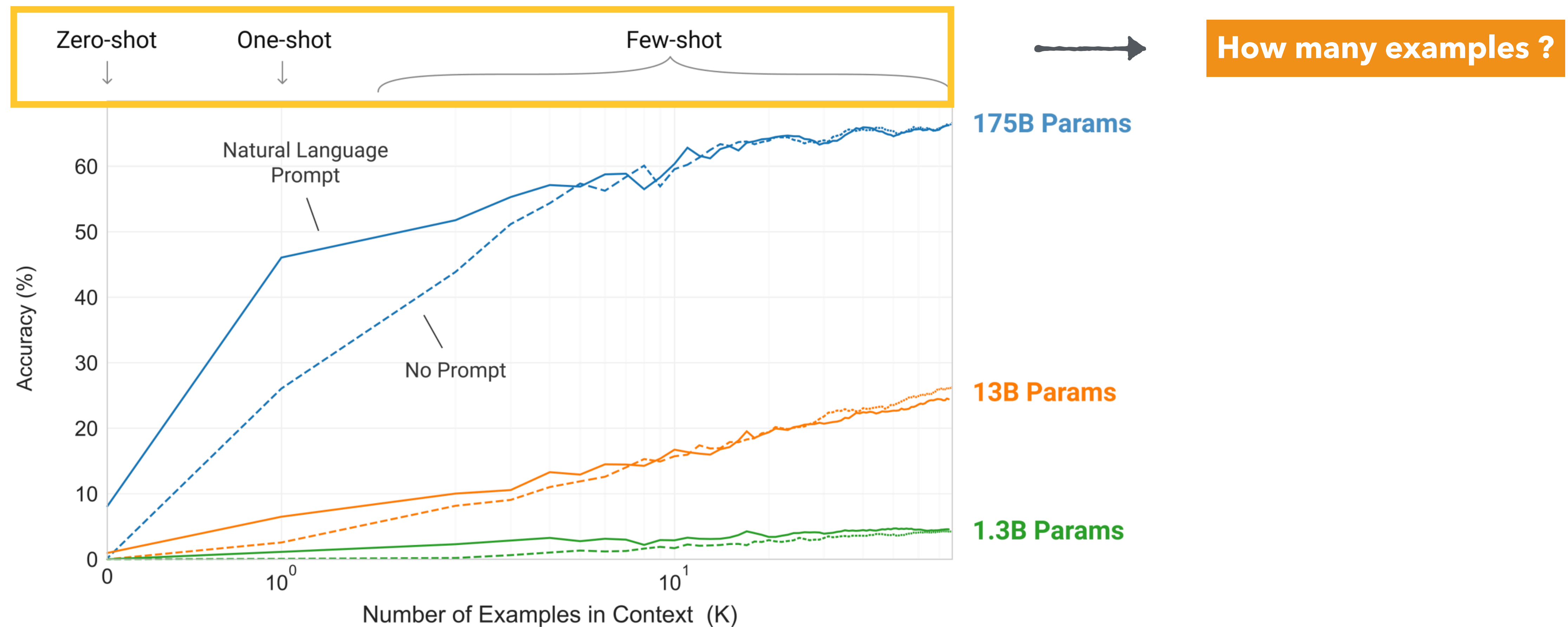
P\* tuning

LM + P\* tuning

# What is in-context learning?

The input to the model describes a new task with some possible examples, **in natural language**.

Effective on **very large models** (173B GPT-3)



# In-context learning: Prompt types

## Zero-shot

Natural language descriptions only

- 1 Find the nationality of people: — *Task description*
- 2 Marie Curie => — *Task*


## One-shot

Description + one example

- 1 Find the nationality of people: — *Task description*
- 2 Albert Einstein => German — *Example*
- 3 Marie Curie => — *Task*

## Few-shot

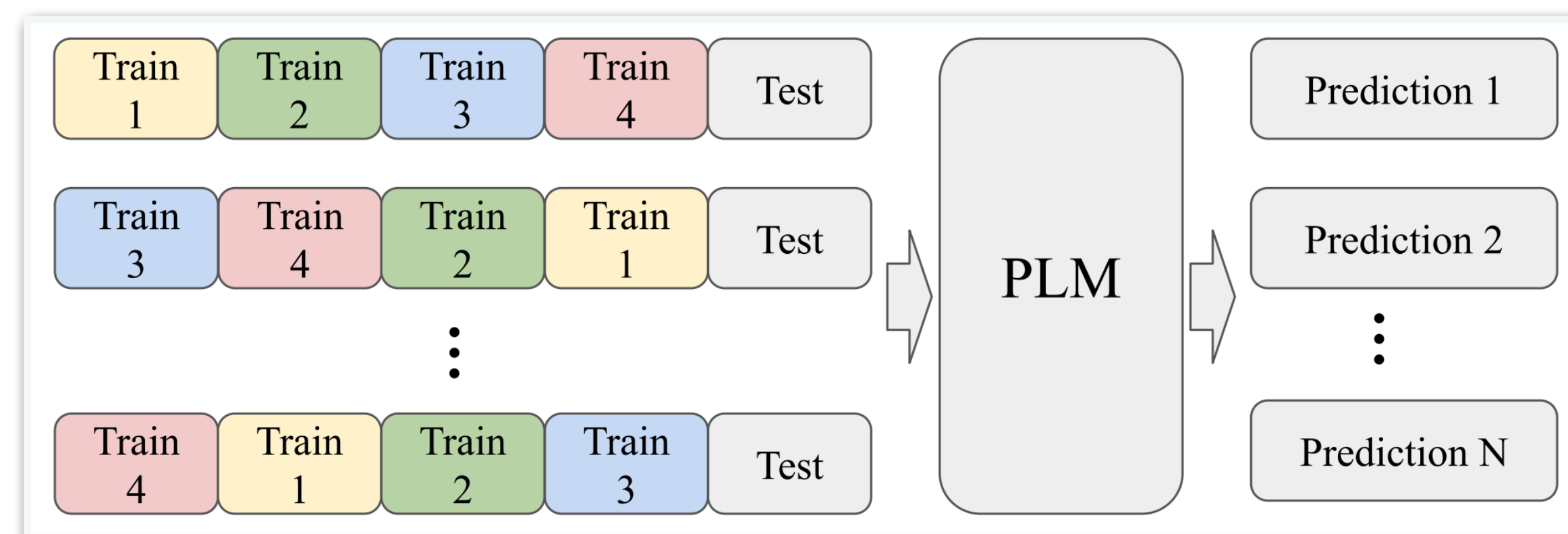
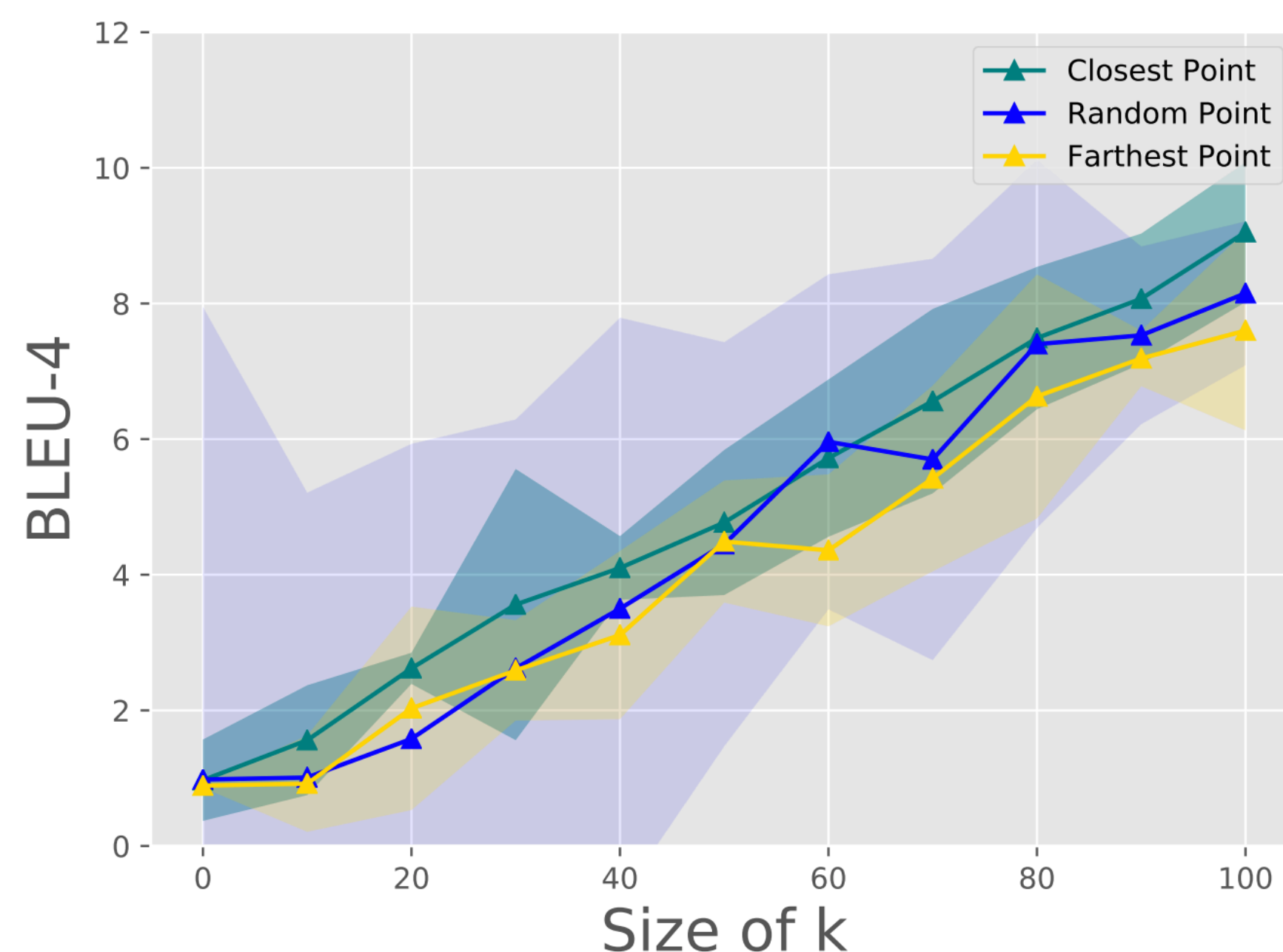
Description + a few example (3-100)  
*[5-10 is most common]*

- 1 Find the nationality of people: — *Task description*
  - 2 Albert Einstein => German — *Examples*
  - 3 Alan Turing => English
  - 4 Mahatma Gandhi => Indian
  - 5 Marie Curie => — *Task*
- 

**How to make ?**

# Challenge: which sets of examples?

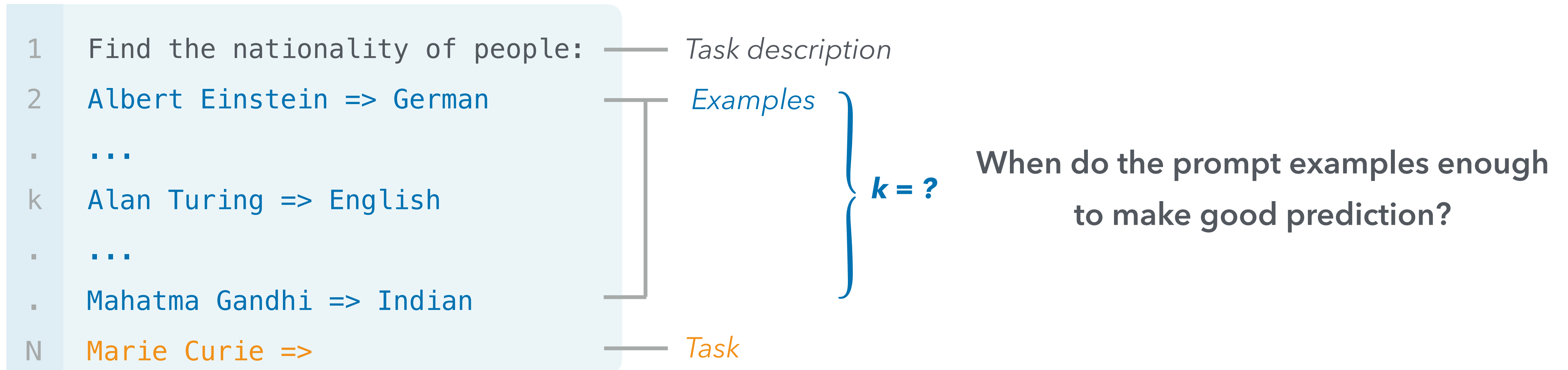
Let's assume users are given a training data set to choose prompt examples.



**Different** few-shot **example sets** lead to very different results.

**Different ordering** of the same set also lead to different results!!

# Challenge: when "enough" examples?



The model **performs better** when the **test input** is **similar** to some **training input**.  
But it's **hard to get coverage** in 30 examples.

# Research objectives

We present, **ScatterShot**, to help users *interactively and iteratively find high-quality demonstrative examples to build effective in-context functions.*



# Scattershot principles

1 Handle common patterns →

Help the user **discover** previously **unexplored patterns**.

2 Not neglect unusual ones →

Help the user **prioritize** the most **informative** examples.

3 Cost effective →

**Minimize** annotation **cost**.

# User interface

The screenshot shows a user interface for date extraction. At the top, a task description reads: "Extract all the mentioned dates as detailed as possible, in the ISO". Below this is a section for "EXISTING FEW-SHOT EXAMPLES" with a "PREVIEW" button and a "COUNT: 3" indicator. Three examples are shown, each with an original prompt (O) and a predicted prompt (P). The first example shows "today" being converted to "2000-01-05". The second shows "Oct. 23, 1999" being converted to "1999-10-23". The third shows "N/A" for the prediction. A red box labeled "Detected data phrases" points to the date parts in the prompts. A green box labeled "LLM generations" points to the predicted dates. Below this is a "CANDIDATES" section with a "LOAD A NEW BATCH" button and a "REMAINING INSPECTION BUDGET: 200" indicator. Two candidate examples are shown. The first has a predicted date of "2000-11-25" and is labeled "Good examples". The second has a predicted date of "1996" and is labeled "Bad examples". On the right side, arrows point from the task description, the examples section, and the candidate batches to their respective labels. The "Good examples" and "Bad examples" labels are highlighted in yellow.

Extract all the mentioned dates as detailed as possible, in the ISO

EXISTING FEW-SHOT EXAMPLES PREVIEW **Detected data phrases** COUNT: 3

O [Posted: 2000-01-05] Photo: today .  
P today == 2000-01-05

O [Posted: 1989-10-31] Slepian was killed on Oct. 23, 1999 .  
P Oct. 23, 1999 == 1999-10-23

O [Posted: 1989-10-31] It hopes to contribute to business  
P N/A

LLM generations

CANDIDATES LOAD A NEW BATCH REMAINING INSPECTION BUDGET: 200

O [Posted: 2000-01-06] He was plucked on Thanksgiving Day  
P Thanksgiving Day == 2000-11-25

O [Posted: 1998-02-27] nineteen ninety-six in Atlanta.  
P nineteen ninety-six == 1996

Task description

Prompt examples

Candidate batches

Good examples

Bad examples

How can we use the **least examples** to cover **most prompt patterns**?

# Scattershot algorithm

Input-output pairs, iteration 1 to  $i - 1$

[Posted: 1998-02-27] **nineteen ninety-six** in Atlanta  
nineteen ninety-six == 1996

[Posted: 2000-01-05] **Photo:** on **today** .  
today == 2000-01-05

[Posted: 2000-01-06] **He was plucked on Thanksgiving** Day.  
Thanksgiving == 1999-11-25

**A** Existing prompt examples

# Scattershot algorithm

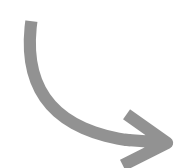
Input-output pairs, iteration 1 to  $i - 1$

[Posted: 1998-02-27] **nineteen ninety-six** in Atlanta  
nineteen ninety-six == 1996

[Posted: 2000-01-05] **Photo:** on **today** .  
today == 2000-01-05

A

[Posted: 2000-01-06] **He was plucked on Thanksgiving** Day.  
Thanksgiving == 1999-11-25



Key phrase templates

**PRON** (Halloween, Thanksgiving)

**DATE** (today, Oct. 23, 1999)

**NUM years ago** (24 years ago)

B

Extract key phrases & slices

# Slice-based Sampling

## Input-output pairs, iteration 1 to $i - 1$

[Posted: 1998-02-27] **nineteen ninety-six** in Atlanta  
nineteen ninety-six == 1996

[Posted: 2000-01-05] **Photo:** on **today** .  
today == 2000-01-05

[Posted: 2000-01-06] **He was plucked on Thanksgiving Day**.  
Thanksgiving == 1999-11-25

### Key phrase templates

**PRON** (Halloween, Thanksgiving)

**DATE** (today, Oct. 23, 1999)

**NUM years ago** (24 years ago)

## B Extract key phrases & slices

## Key phrases & data slices, iteration $i$

C

- ① ✓ [Posted: 1998-02-27] Atlanta nineteen ninety-six.  
X [Posted: 1989-10-31] It hopes to control 5% of jewelry business  
? [Posted: 2013-10-02] 19 - 20 October, Chevron House.
- ② ? [Posted: 2014-12-25] @viereedom Merry Christmas!  
? [Posted: 2014-10-12] HALLOWEEN SHOW FOR HSBC FAMILY...  
X [Posted: 2000-01-06] He was plucked on Thanksgiving Day.
- ③ ? [Posted: 2015-03-21] Her last run was 24 years ago  
✓ [Posted: 2014-07-09] Photo: One year ago, #Singapore  
X [Posted: 2015-04-20] But it's already 10 months ago!!
- ④ ✓ [Posted: 2015-01-02] Are you going to yoga today?  
? [Posted: 2000-01-05] Photo: today.  
✓ [Posted: 2014-10-19] Lunch at Agnes B Cafe yesterday.

# Prioritize sampled examples

Input-output pairs, iteration 1 to  $i - 1$

[Posted: 1998-02-27] **nineteen ninety-six** in Atlanta  
 nineteen ninety-six == 1996

[Posted: 2000-01-05] **Photo:** on **today** .  
 today == 2000-01-05

[Posted: 2000-01-06] **He was plucked on Thanksgiving Day.**  
 Thanksgiving == 1999-11-25

**A**

Key phrase templates

**PRON** (Halloween, Thanksgiving)  
**DATE** (today, Oct. 23, 1999)  
**NUM years ago** (24 years ago)

**B** Extract key phrases & slices

Key phrases & data slices, iteration  $i$

**C**

①	✓	[Posted: 1998-02-27] Atlanta nineteen ninety-six.	$n=449$
	X	[Posted: 1989-10-31] It hopes to control 5% of jewelry business	$m=10$
	?	[Posted: 2013-10-02] 19 - 20 October, Chevron House.	$k=4$
			$\mu=4.82$
②	?	[Posted: 2014-12-25] @viereedom Merry Christmas!	$n=19$
	?	[Posted: 2014-10-12] HALLOWEEN SHOW FOR HSBC FAMILY...	$m=2$
	X	[Posted: 2000-01-06] He was plucked on Thanksgiving Day.	$k=0$
			$\mu=4.34$
③	?	[Posted: 2015-03-21] Her last run was 24 years ago	$n=31$
	✓	[Posted: 2014-07-09] Photo: One year ago, #Singapore	$m=5$
	X	[Posted: 2015-04-20] But it's already 10 months ago!!	$k=1$
			$\mu=3.61$
④	✓	[Posted: 2015-01-02] Are you going to yoga today?	$n=113$
	?	[Posted: 2000-01-05] Photo: today.	$m=3$
	✓	[Posted: 2014-10-19] Lunch at Agnes B Cafe yesterday.	$k=3$
			$\mu=1.14$

**Prioritize** similar data that has **low performance**, are **large**, and slices that have **not been** sampled many times.

$$\mu_{i,c} = \underbrace{\left(1 - \frac{k}{m}\right)}_{\text{Error rate}} \cdot \underbrace{\ln n}_{\text{Size}} + \underbrace{\sqrt{\frac{\ln t}{m}}}_{\text{Sample Rarity}}$$

Slice  $c$  has  $n$  examples,  $m$  are labeled in previous iterations. Out of  $m$ , the current function is correct on  $k$ .

# How to handle no ground truth labels?

We estimate function quality by re-ordering stability.

[Posted: 2014-12-25] @viereedom Merry Christmas! A

 Unanimity voting


- Christmas == 2014-12-25
- Christmas == 2014-12-25
- Christmas == 2014-12-25

 Manual inspection


Keep Christmas == 2014-12-25

*Annotations: A blue checkmark and arrow point to the top item. A blue arrow points down from the top item to the manual inspection section. A blue 'X' is at the bottom of the manual inspection section.*

[Posted: 1998-02-27] Atlanta nineteen ninety-six. B

 Unanimity voting

- nineteen ninety-six == 1996-01
- nineteen ninety-six == 1996
- 1996 == 1996

 Manual inspection

Edit nineteen ninety-six == 1996

*Annotations: A blue checkmark and arrow point to the top item. A blue arrow points down from the top item to the manual inspection section. A blue 'X' is at the bottom of the manual inspection section. The text '1996' in the manual edit is highlighted in green.*

# Scattershot evaluation

## Task & Datasets

### 1 Simulation Experiment

- Simulate the labeling process

### 2 Within-subject User Study

- 10 person evaluation
- QA-pair rewriting task

### Temporal Expression Extraction

O [Posted: 2000-01-05] Photo: today .

P today == 2000-01-05

O [Posted: 1989-10-31] Slepian was killed on Oct. 23, 1999 .

P Oct. 23, 1999 == 1999-10-23

O [Posted: 1989-10-31] It hopes to control 5% of jewelry business

P N/A

### Question-Answer Pair Rewriting

O Q: Where are the buildings? A: in distance

P Q: Are the buildings in distance? A: yes

O Q: Why is it dark? A: twilight

P Q: Is it dark because of the twilight? A: yes

O Q: Is the water warm or cold? A: cold

P Q: Is the water cold? A: yes



# 1 Simulation performance

## Temporal

Conditions	Extraction			Normalization		
	F1	Precision	Recall	F1	Precision	Recall
Random	73.2 ± 4.0	74.0 ± 3.8	72.9 ± 4.1	66.8 ± 3.2	67.3 ± 3.3	67.0 ± 3.1
SCATTERSHOT	75.0 ± 2.9	75.6 ± 2.8	74.7 ± 2.9	70.9 ± 3.4**	71.3 ± 3.5*	71.2 ± 3.2**

## QA-Pair

Conditions	ROUGE-L	BLEU-4
Rule-based	78.4	66.7
Random	74.3 ± 3.9	65.4 ± 3.5
SCATTERSHOT	80.0 ± 3.5*	69.1 ± 3.1*

The significant improvements, measured by the student's **t-test** are marked with **\***:  $p < 0.05$ , and **\*\***:  $p < 0.01$ .

### Quantitative Results:

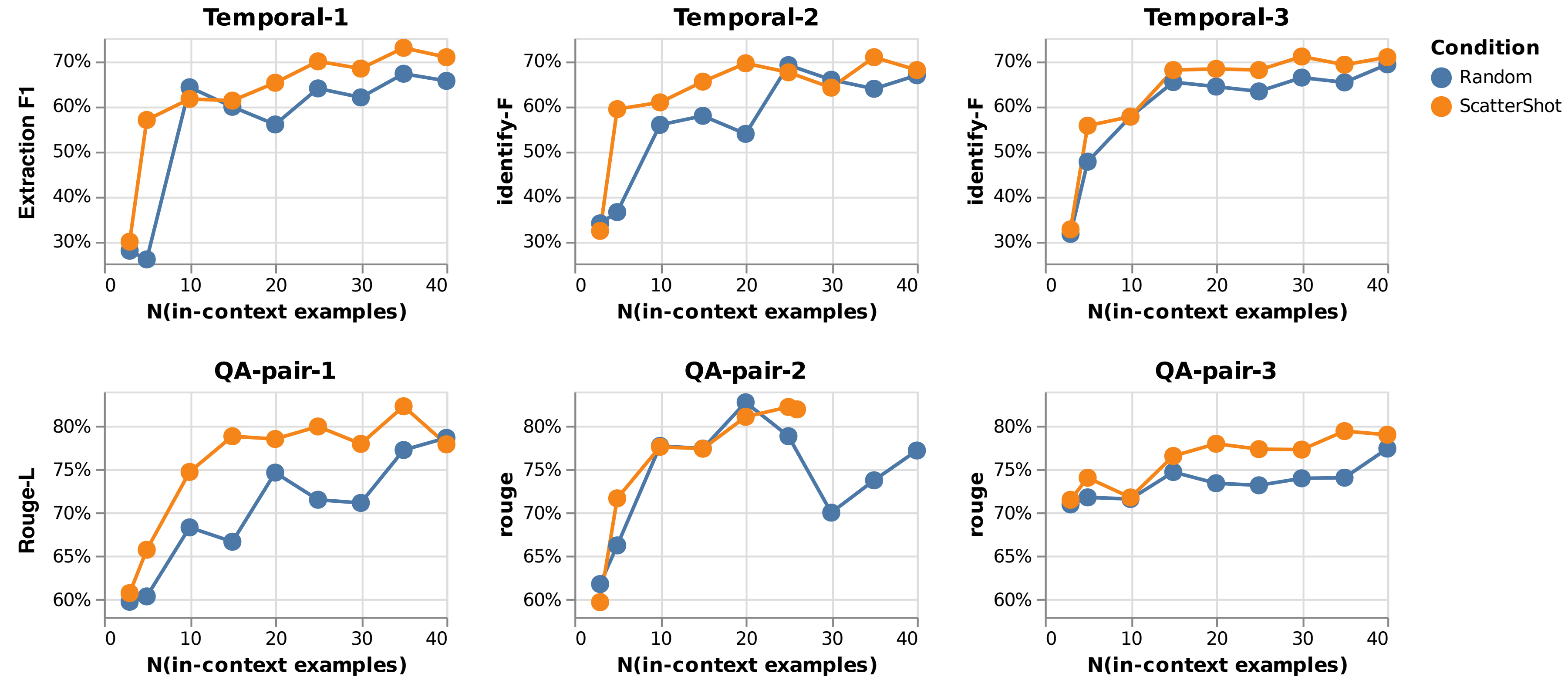
✓ Compared with the Random condition, **ScatterShot outperformed** the baseline on **all metrics**.

# Example outputs

<b>Coverage:</b> Transforms more forms of inputs.		
Input	Q: Are there more girls or boys? A: equal	Q: How many hairs does the sheep in front have? A: infinite
Rule-based	✗ (No generation)	✗ (No generation)
Random	✓ Q: Are the girls and boys equal in number? A: yes	✗ N/A
SCATTERSHOT	✗ Q: Are the girls and boys equal? A: yes	✓ Q: Does the sheep in front have infinite hairs? A: yes
<b>Fluency:</b> Generates outputs that sound natural.		
Input	Q: What make is the phone? A: vtech	Q: What does the woman have on her face? A: headband
Rule-based	✗ Q: Make is the phone vtech? A: yes	✗ Q: Does the woman have on her face headband? A: yes
Random	✓ Q: Is the phone a vtech? A: yes	✓ Q: Does the woman have a headband on her face? A: yes
SCATTERSHOT	✓ Q: Is the phone a vtech? A: yes	✓ Q: Does the woman have a headband on her face? A: yes
<b>Correctness:</b> Produces desired outputs (the new question-pair are logically equivalent to the original pair).		
Input	Q: What monument are they next to? A: unknown	Q: What type of motorcycle is in the picture? A: mountain
Rule-based	✗ Q: Are they next to unknown? A: yes	✗ Q: Is the mountain in the picture ? A: yes
Random	✗ Q: Is the monument unknown? A: yes	✗ Q: Is the mountain type of motorcycle in the picture ? A: yes
SCATTERSHOT	✓ Q: Are they next to an unknown monument? A: yes	✓ Q: Is the motorcycle in the picture a mountain bike? A: yes

Compared with the **Random** condition, and a **Rule-based** system:  
**ScatterShot** functions tend to have better **coverage**, **fluency**, and **correctness**.

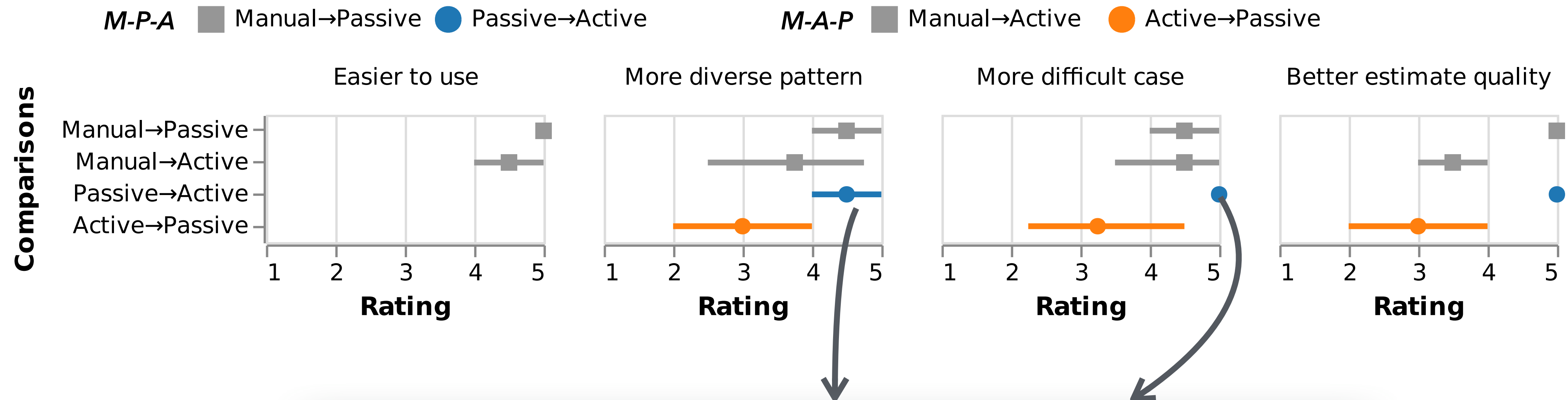
# Performance trajectory w.r.t. examples



We evaluate the **held-out test set** every time we add five more examples to the in-context bucket until the stop condition is satisfied.

**ScatterShot** tends to frequently **outperform** Random, and tends to **have better performance**

## 2 User Study Performance



### **Active learning is effective for humans (More holistic view)!**

*I went through several rounds of pretty similar examples in Step 2 (Random), thinking the function is behaving quite decently, and didn't realize the function needed more diverse and edge cases until I reached Step 3.*

# Performance of user created function

Condition	Step 1	Step 2	Step 3
<i>M-R-S</i>	/ (59.3)	+17.4 (74.7)	<b>+3.2</b> (77.8)
<i>M-S-R</i>	/ (61.8)	<b>+18.1</b> (75.4)	-0.4 (74.9)

(a) ROUGE-L

**R -> S**  
**S -> R**

Condition	Step 1	→ Step 2	→ Step 3
<i>M-R-S</i>	/ (63.9)	<b>+10.1</b> (74.0)	<b>+3.1</b> (76.9)
<i>M-S-R</i>	/ (65.3)	+8.9 (74.2)	-0.6 (73.6)

(b) BLEU-4

**+/-** : represents the **average performance change** compared to the prior step, (number) are the absolute performance.

**M-R-S**: users build in-context functions using methods of "Manual - Random - ScatterShot" in sequence.

**M-S-R**: users use "Manual - ScatterShot - Random" methods in sequence.

***M-R-S** users were able to keep **adding useful examples**, whereas **M-S-R** users **decreased** the function performance by 0.6 in Step 3 (ScatterShot -> Random), indicating that these efforts were wasted.*

# What's more?

- ✓ Slice-based sampling can increase **data space coverage**
- ✗ Random sampling performs less

✓ Interacting with the latest function for users is essential for in-context learning.

✓ Human-AI collaborative labeling for building better functions results in better quality and better task definition.

# Takeaways

**ScatterShot** helps users find *informative input examples* in the unlabeled data, **improves** the *annotator's awareness and handling of diverse patterns*, and ultimately, the *in-context function performance*.

The full user study instructions, and the detailed exit survey, are at:

 **GitHub:** <https://github.com/tongshuangwu/scattershot>

# Thank You!



Sherry

Tongshuang Wu

 [sherryw@cs.cmu.edu](mailto:sherryw@cs.cmu.edu)

 [@tongshuangwu](https://twitter.com/tongshuangwu)



**Hua Shen**

 [huashen218@psu.edu](mailto:huashen218@psu.edu)

 [@huashen218](https://twitter.com/huashen218)



Daniel S. Weld

 [weld@cs.uw.edu](mailto:weld@cs.uw.edu)



Jeffrey Heer

 [jheer@cs.uw.edu](mailto:jheer@cs.uw.edu)



Marco Tulio Ribeiro

 [marcotcr@microsoft.com](mailto:marcotcr@microsoft.com)