

# Exploring Distributional Shifts in Large Language Models for Code Analysis

Shushan Arakelyan<sup>1</sup>, Rocktim Jyoti Das<sup>2</sup>, Yi Mao<sup>3</sup>, Xiang Ren<sup>1</sup>

<sup>1</sup>University of Southern California, <sup>2</sup>IIT Delhi, <sup>3</sup>Microsoft Azure AI



# Motivation

## 20 years ago...

“What characteristics differ between projects used for building predictors?”

Open source	Yes/No
Global development	Yes/No
Code reviews	Yes/No
Static checkers	Yes/No
...	

## Today

More data



Larger models



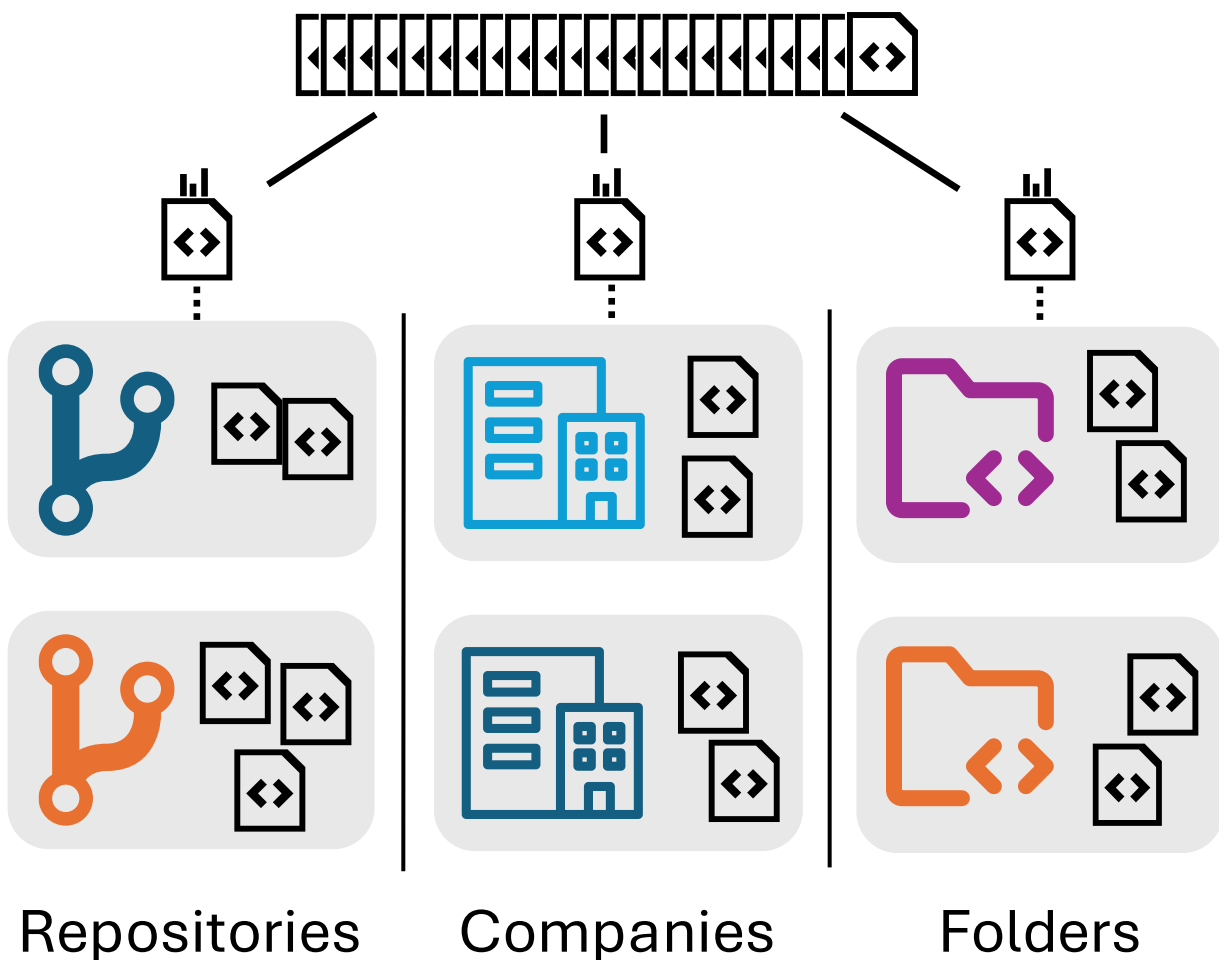
More deployment



Same challenges?

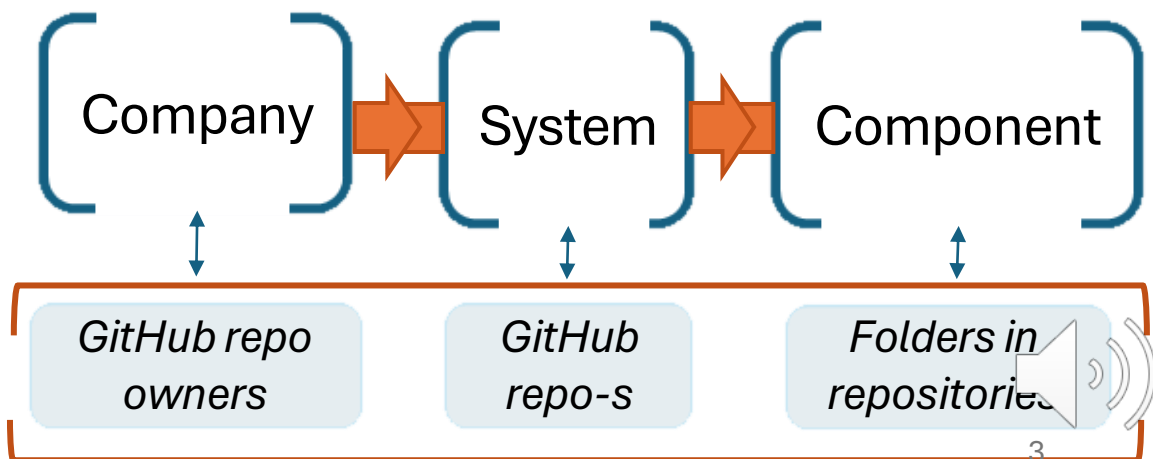


# Data preparation



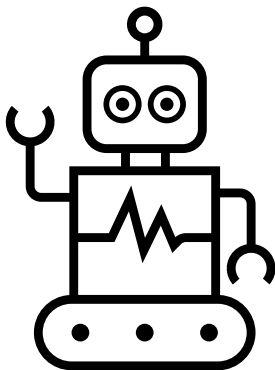
CODE SEARCH	TRAINING	DEV	TESTING
GO	635,635	28,483	14,291
JAVA	908,886	30,655	26,909
<b>JAVASCRIPT</b>	<b>247,773</b>	<b>16,505</b>	<b>6,483</b>
PHP	1,047,406	52,029	28,391
PYTHON	824,342	46,213	22,176
RUBY	97,580	4,417	2,279

Train	
org.	9737
repos.	15858
fold.	25268

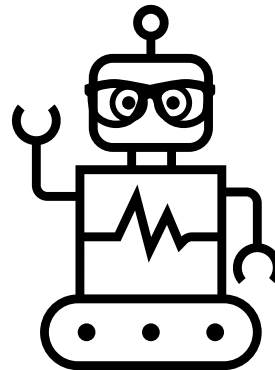


# Experimental setup

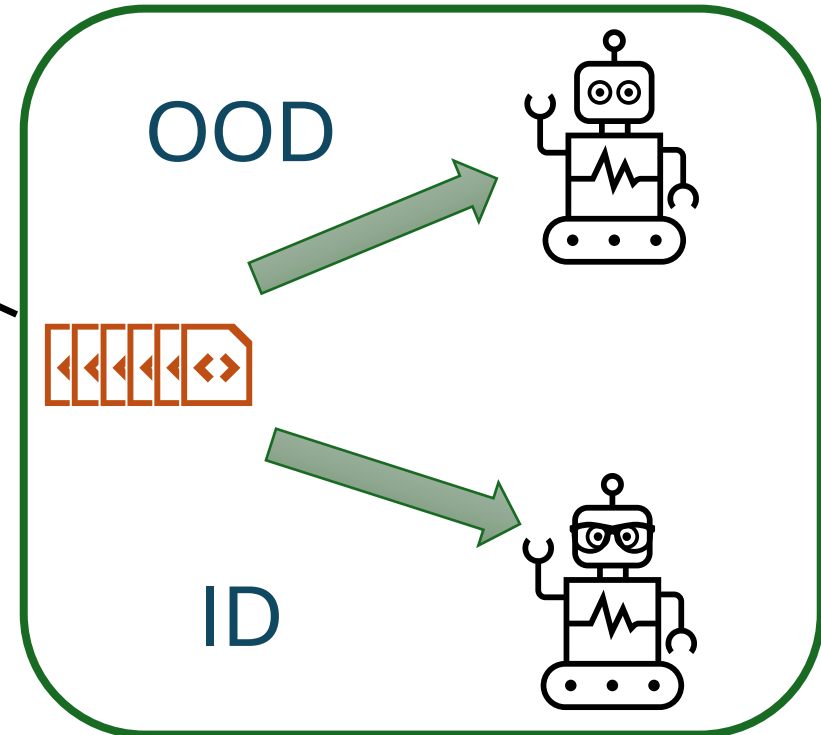
Train data



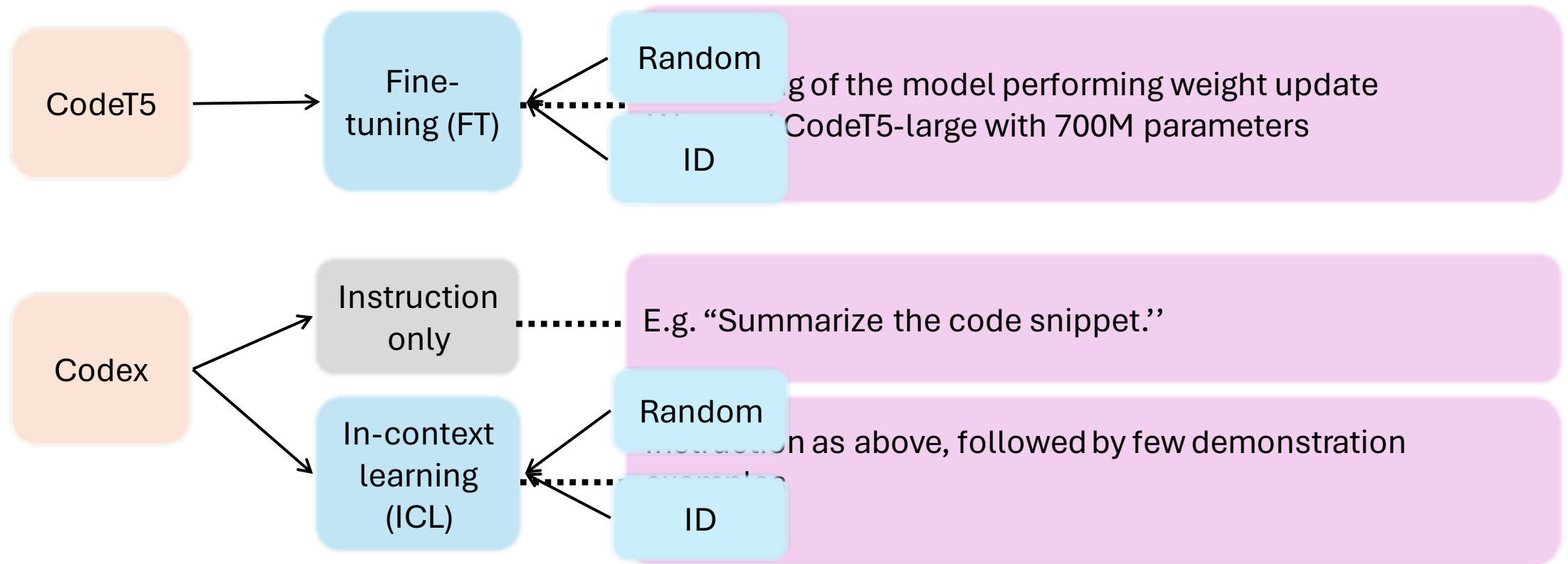
Target domain



Test both models on  
unseen samples from  
the target domain

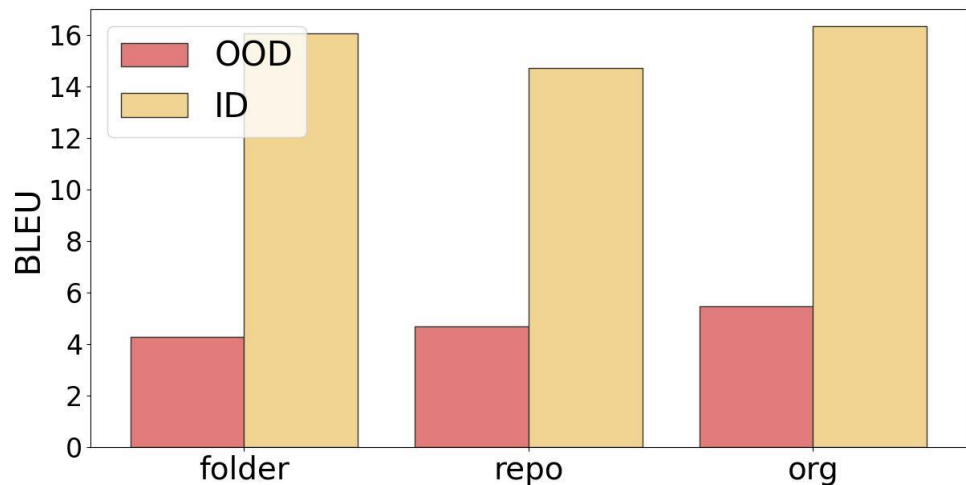


# Models and methods



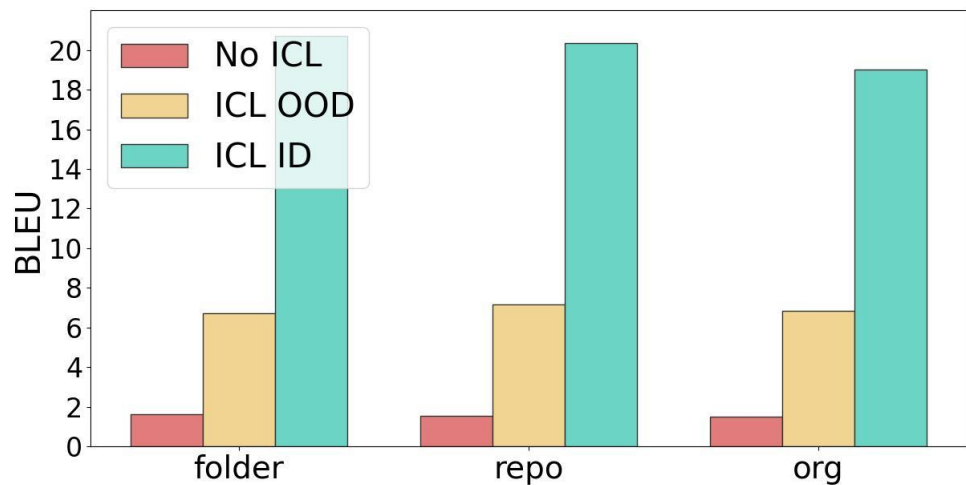
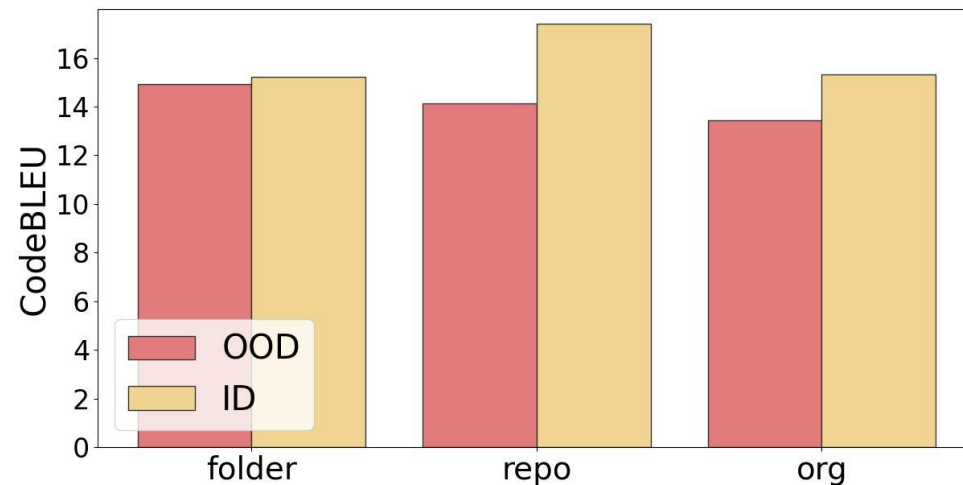
# Results: Performance ID vs OOD

## Code summarization

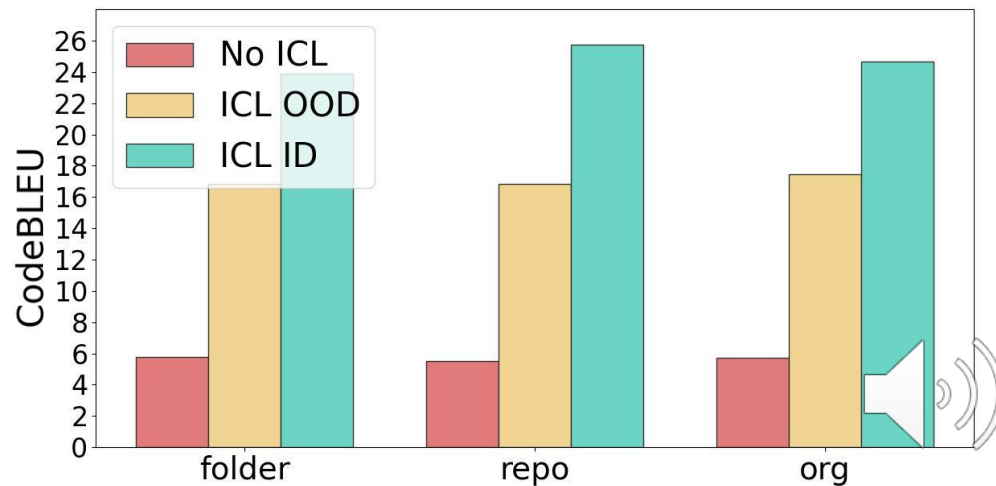


CodeT5

## Code generation

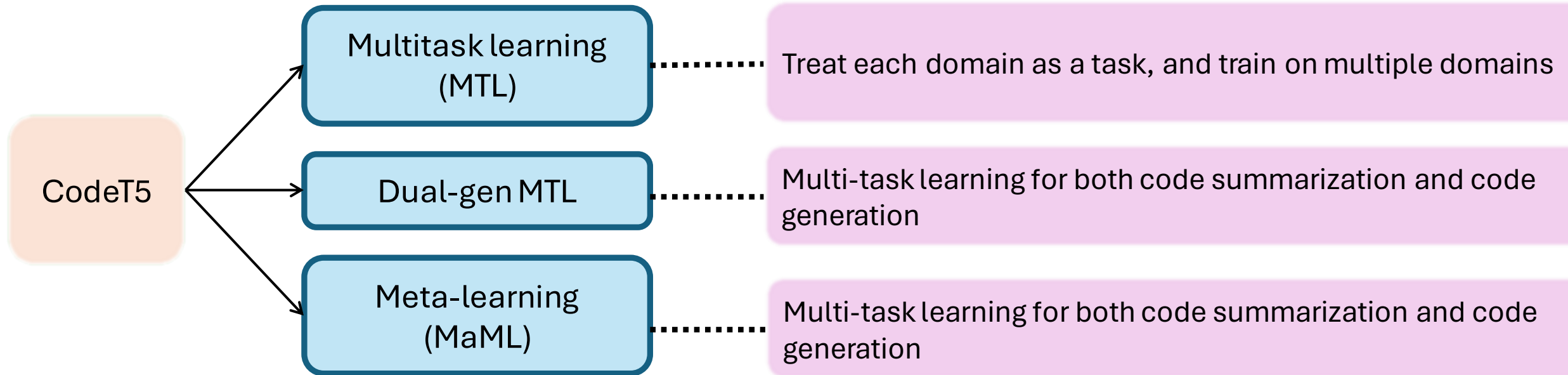


Codex

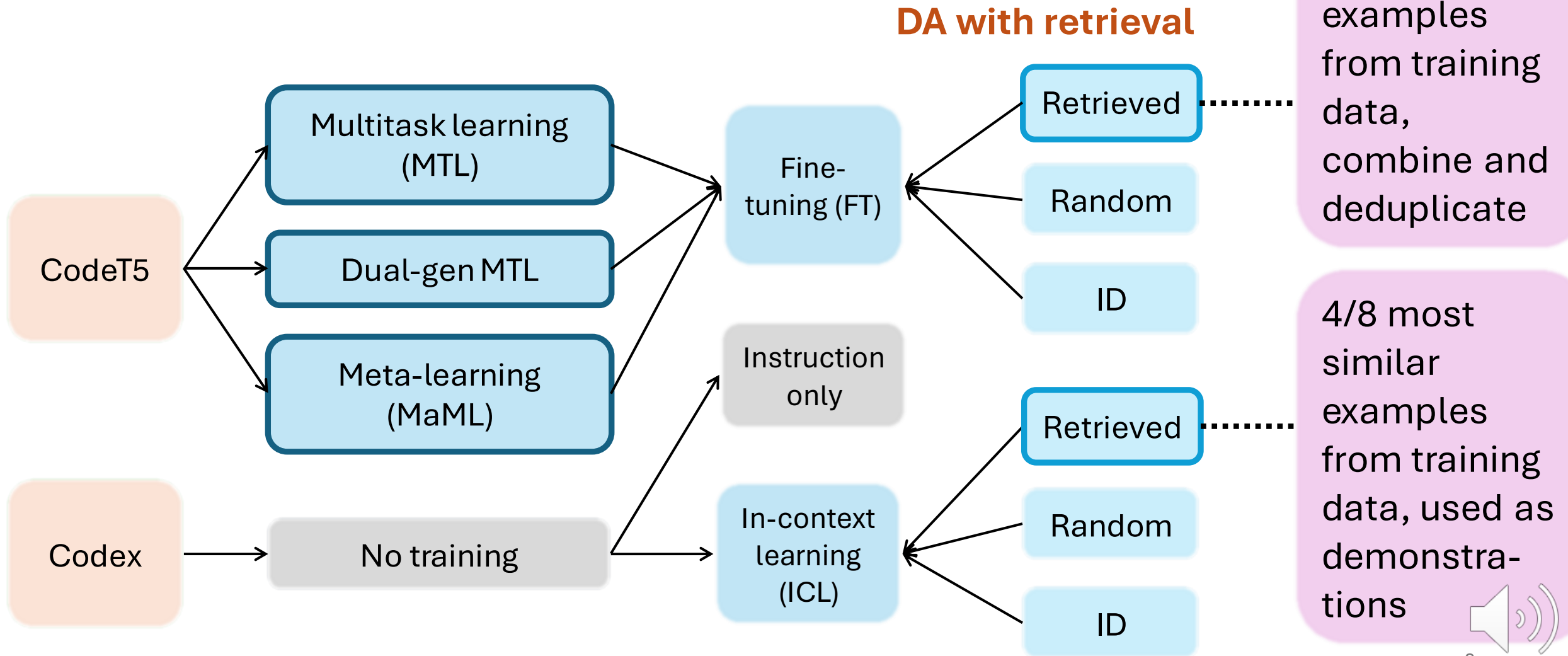


# How to improve OOD performance?

## Training



# How to improve OOD performance?

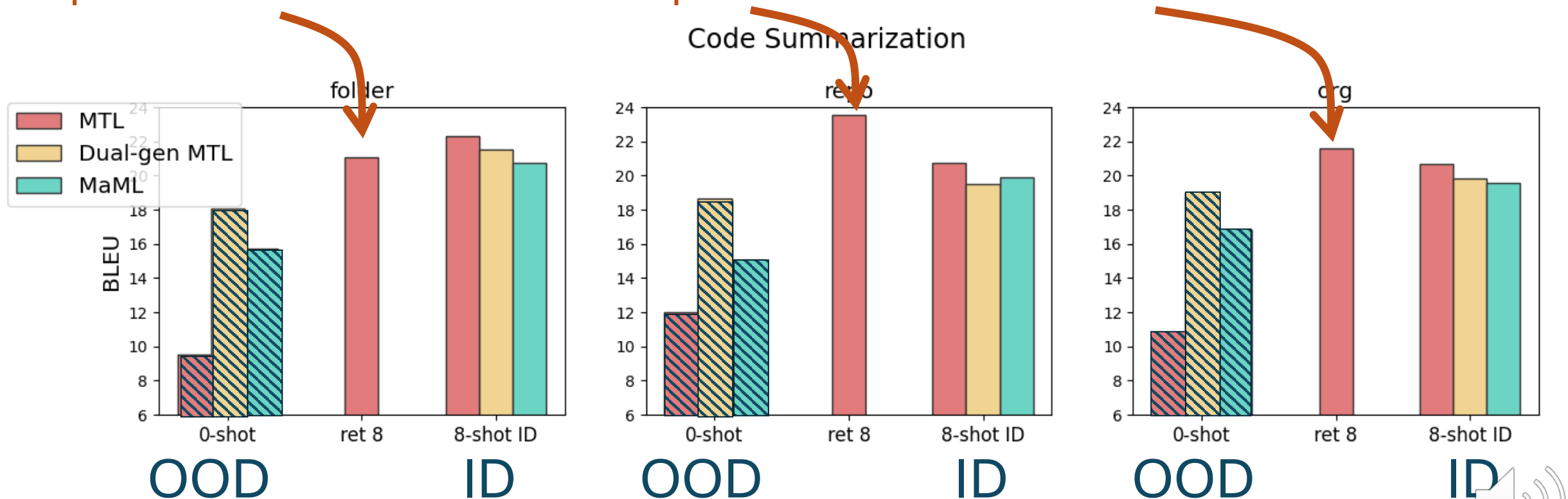




# Results [CodeT5]

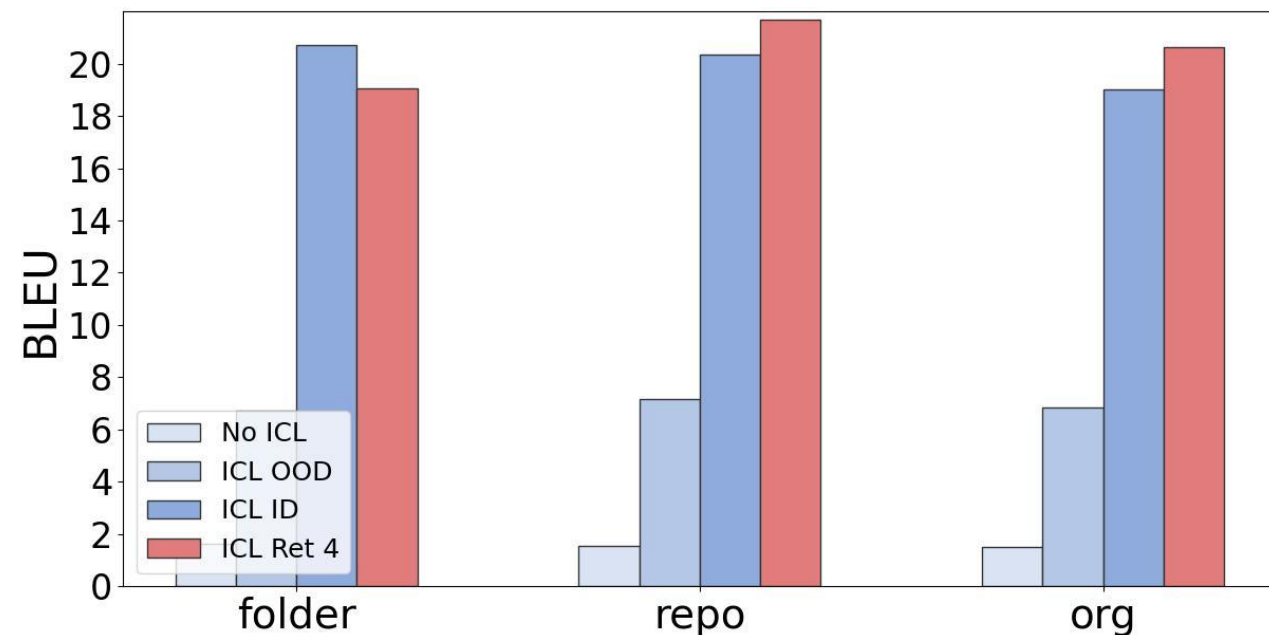
Training does not get rid of ID vs OOD performance discrepancy

Supervision with retrieved examples is more effective!

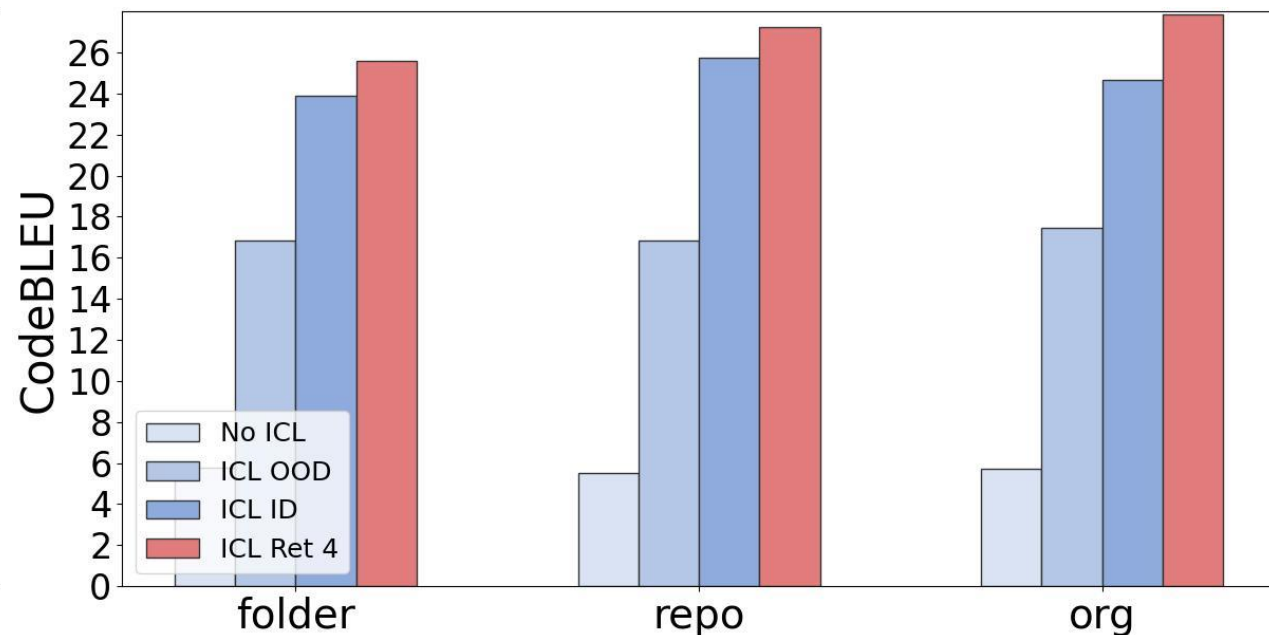


# Results [Codex]

## Code summarization



## Code generation



Supervision with retrieved examples is effective with ICL



# Findings

- Splits naturally occurring in software present **distributional shift challenge**
- Domain adaptation can be effective with a very small amount of data
- Retrieving examples for supervision is effective in **combating distribution shift**

