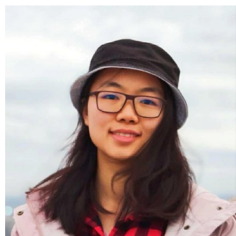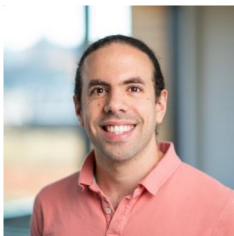# FiD-ICL: A Fusion-in-Decoder Approach for Efficient In-Context Learning

Qinyuan Ye     Iz Beltagy     Matthew E. Peters     Xiang Ren     Hannaneh Hajishirzi
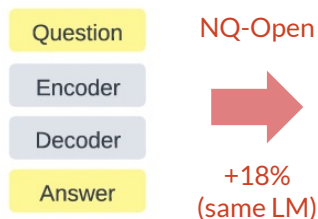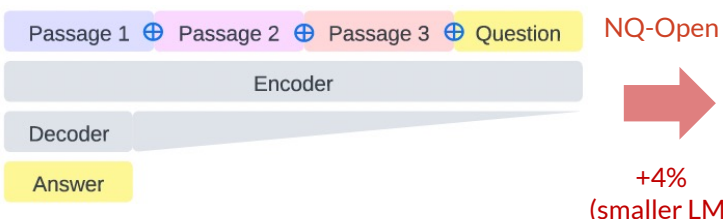
# Background: QA vs. ICL

## Closed-book QA
(Roberts et al., 2020)



NQ-Open

+18%
(same LM)

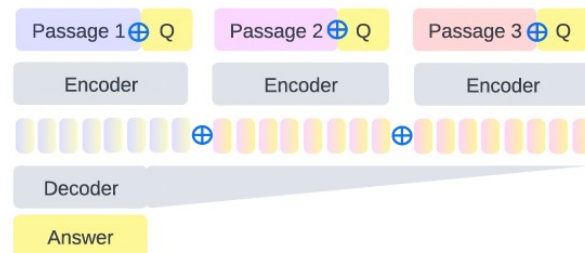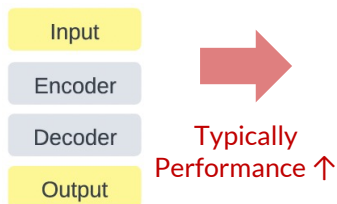## Retrieval-Augmented Generation
(Lewis et al., 2020)



NQ-Open

+4%
(smaller LM)

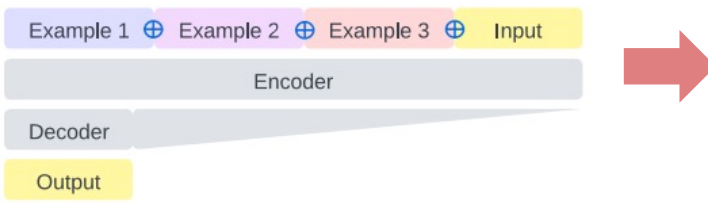## Fusion-in-Decoder
(Izacard et al., 2020)
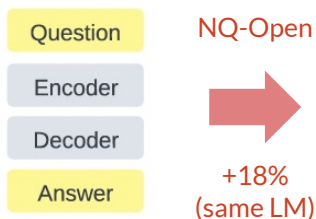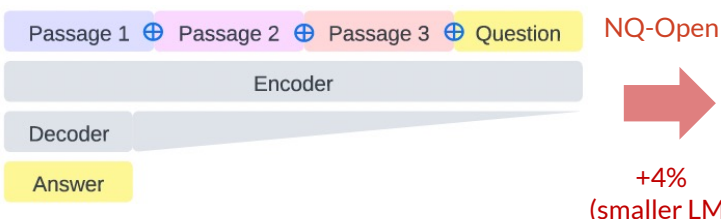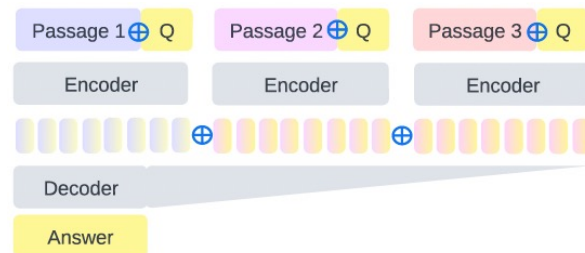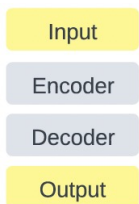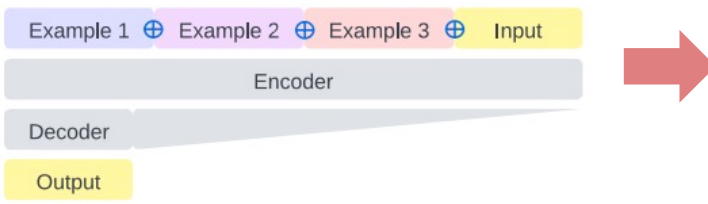


## Zero-shot Learning



Typically
Performance ↑

## Few-shot In-Context Learning



3

# Background: QA vs. ICL

## Closed-book QA
(Roberts et al., 2020)

| Question |
| Encoder |
| Decoder |
| Answer |

NQ-Open → +18% (same LM)

## Retrieval-Augmented Generation
(Lewis et al., 2020)

Passage 1 ⊕ Passage 2 ⊕ Passage 3 ⊕ Question
Encoder
Decoder
Answer

NQ-Open → +4% (smaller LM)

## Fusion-in-Decoder
(Izacard et al., 2020)

Passage 1 ⊕ Q    Passage 2 ⊕ Q    Passage 3 ⊕ Q
Encoder          Encoder          Encoder
Decoder
Answer

## Zero-shot Learning

| Input |
| Encoder |
| Decoder |
| Output |

Typically Performance ↑

## Few-shot In-Context Learning

Example 1 ⊕ Example 2 ⊕ Example 3 ⊕ Input
Encoder
Decoder
Output

## This work: FiD-ICL

Example 1    Example 2    Example 3    Input
Encoder      Encoder      Encoder      Encoder
Decoder
Output

USC University of Southern California

AI2

# FiD-ICL



Fusion-in-decoder

⊕ "Intermediate Fusion"
Concat. Hidden Repr.

# FiD-ICL



can be pre-computed

**Fusion-in-decoder**

⊕ *"Intermediate Fusion"*

Concat. Hidden Repr.

Example 1 | Example 2 | Example 3 | Input

Encoder | Encoder | Encoder | Encoder

Decoder

Output

# Compared Methods



Referred to as "fusion" methods for ICL
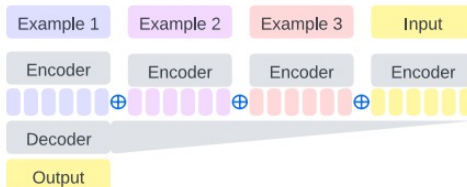
**Concat-based ICL**
⊕ *"Early Fusion"*
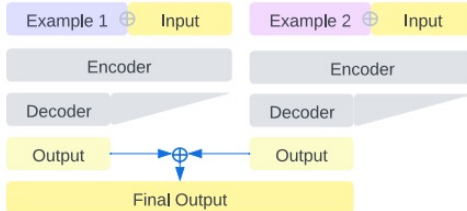Concat. Raw Text

**Fusion-in-decoder**
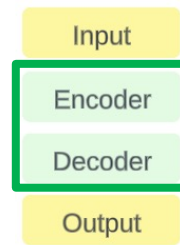⊕ *"Intermediate Fusion"*
Concat. Hidden Repr.

**Ensemble-based ICL**
⊕ *"Late Fusion"*
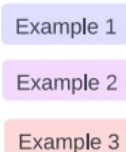Aggregate Scores for Rank Classification

# Compared Methods
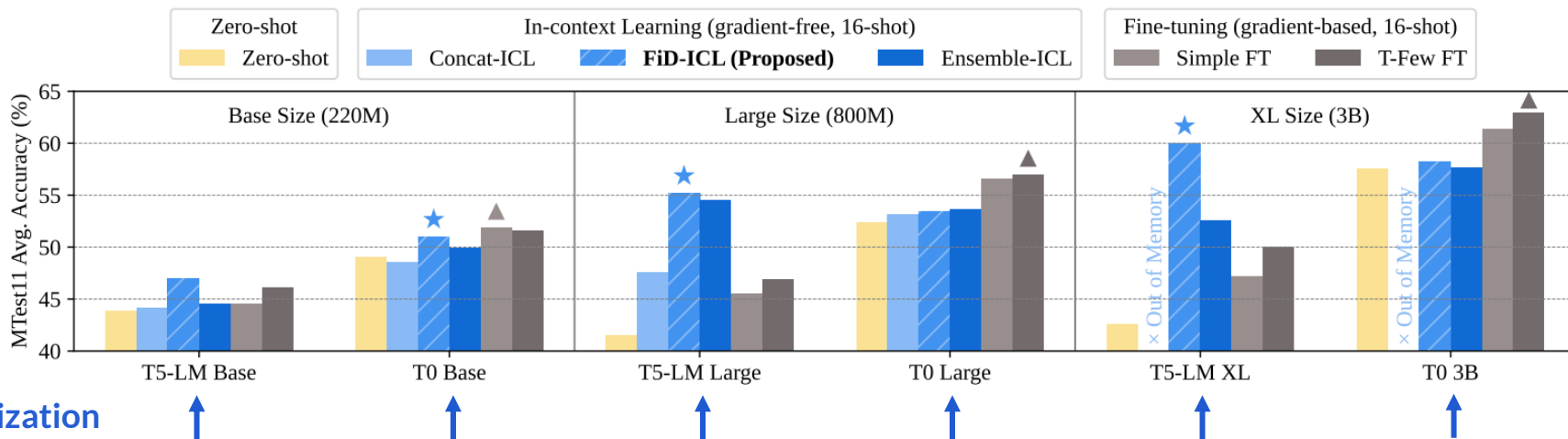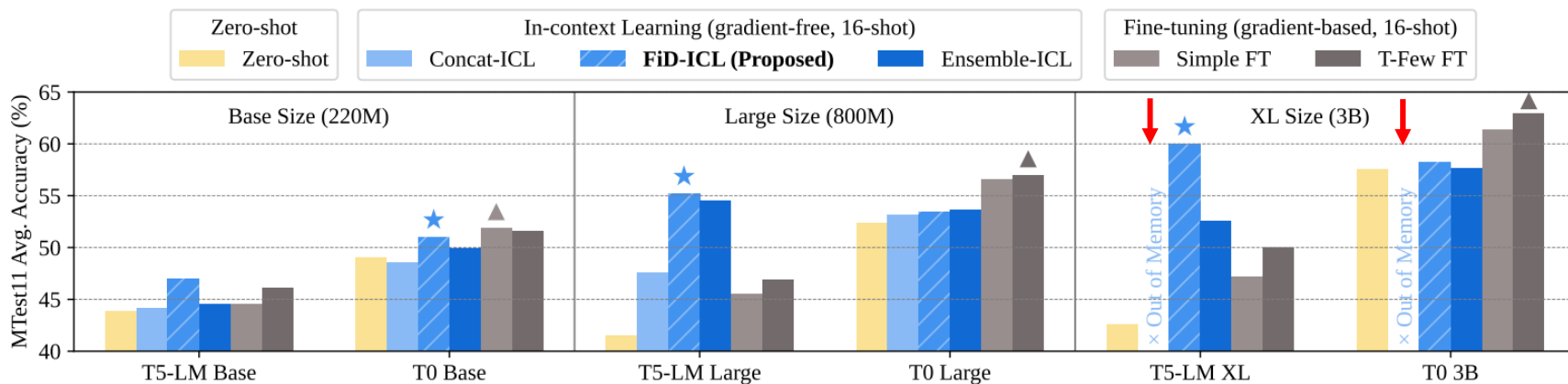
# Main Results



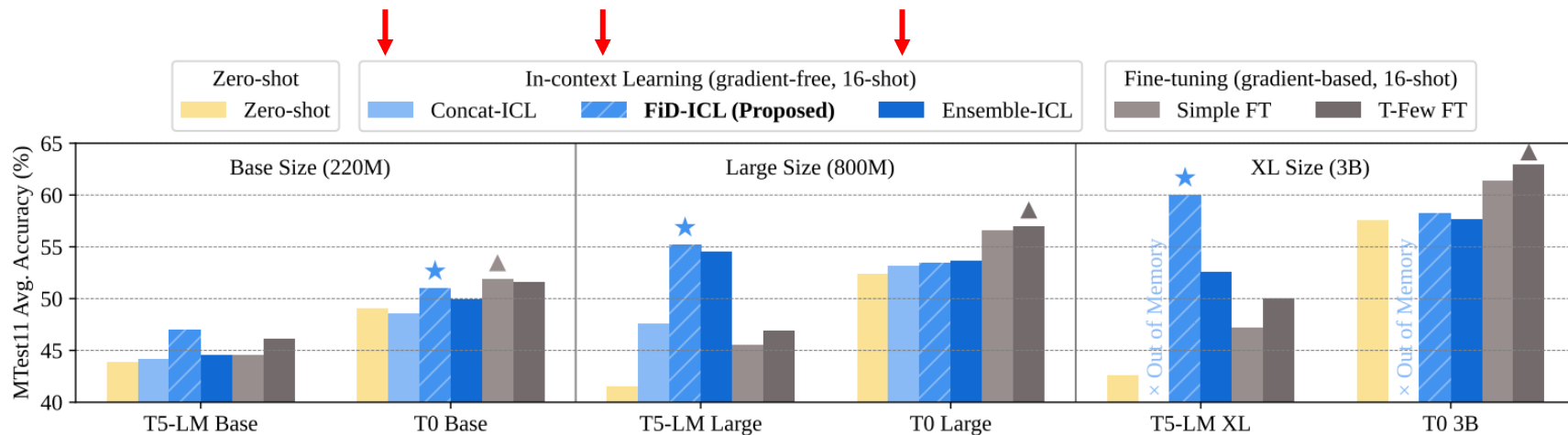Using *Public Pool of Prompts (P3)* dataset
Using a *meta-training* setting
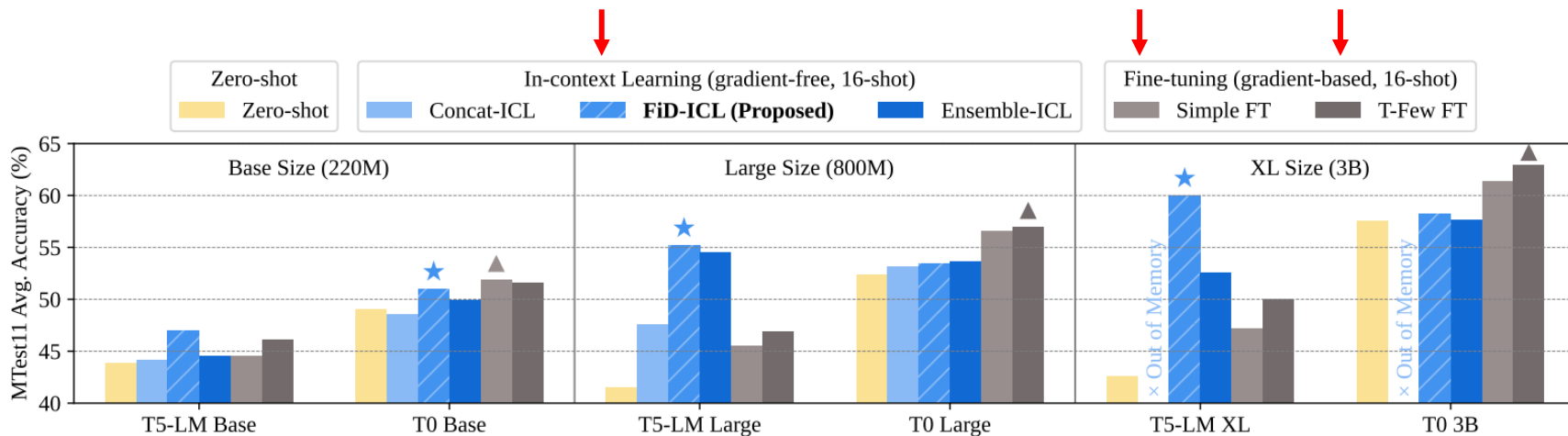
# Main Results



**FiD-ICL enables efficient meta-training (Concat-ICL would fail at 3B)**
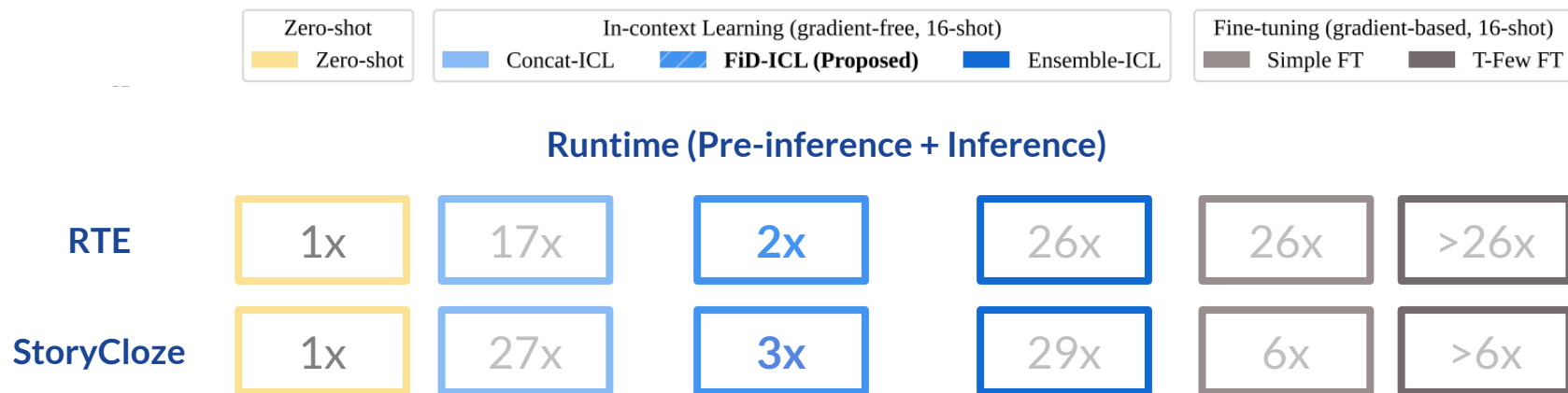
# Main Results



FiD-ICL outperforms the other two fusion methods (Concat and Ensemble)

# Main Results



The gap between FiD-ICL (★ gradient-free) and fine-tuning (▲ gradient-based) is *<3%*.

# Efficiency

*Limitations apply. Fine-tuned models are still more efficient for large-scale inference.*

# Analysis (or... surprise?)
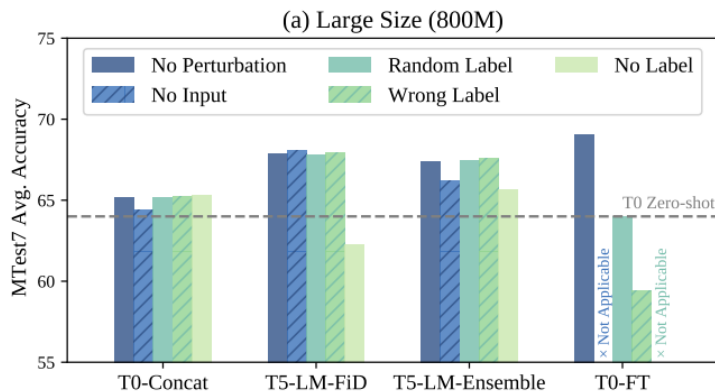
## Number of Shots



Average performance **does not** grow with more shots.

It's **task-dependent**.

## Perturbation to In-context Examples

(Inspired by Min et al., 2022)



Performance is rather **insensitive** to perturbations to in-context examples.

Still **not** learning effectively.

# Conclusion

**FiD-ICL**, a fusion-in-decoder approach for efficient in-context learning

**Performance**
It outperforms Concat-ICL and Ensemble-ICL.
The gap between FiD-ICL and fine-tuning is **<3%** on P3 meta-test tasks.

**Efficiency**
FiD-ICL is more efficient than Concat-ICL, Ensemble-ICL.
More efficient than fine-tuning when considering pre-inference + inference time*.

**Limitations**
FiD-ICL is still *not perfect*; still has the similar limitations as Concat-ICL.

**Implications**
Insights and methodologies from *open-domain QA* are very useful!
FiD-ICL is related to *retrieval augmentation*, *sparse attention*, and *hypernetworks*.

USC University of Southern California

Ai2