# Investigating the Effectiveness of Multiple Expert Models Collaboration

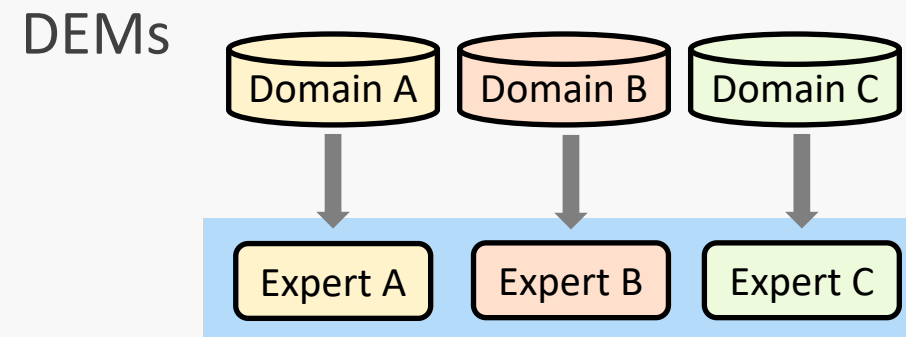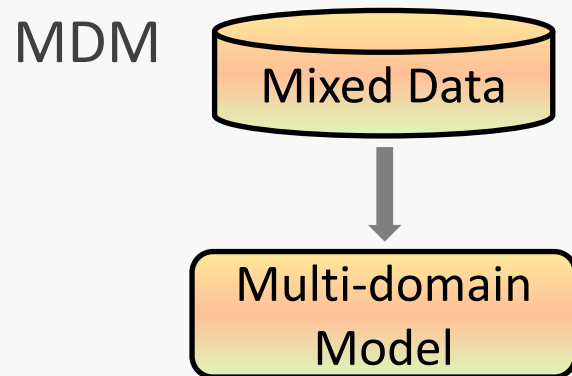**Ikumi Ito[1], Takumi Ito[1, 2], Jun Suzuki[1, 3], Kentaro Inui[4, 1, 3]**

[1]Tohoku university, [2]Langsmith Inc., [3]RIKEN, [4]MBZUAI

# Explore the potential of multiple models in multi-domain translation

- One challenge in machine translation is multi-domain adaptation [Saunders+ '22]

- Main approaches for multi-domain
    - a single Multi-Domain Model (MDM)
    - multiple Domain Expert Models (DEMs)

- The contributions of this work
    - Demonstrated effectiveness of DEMs
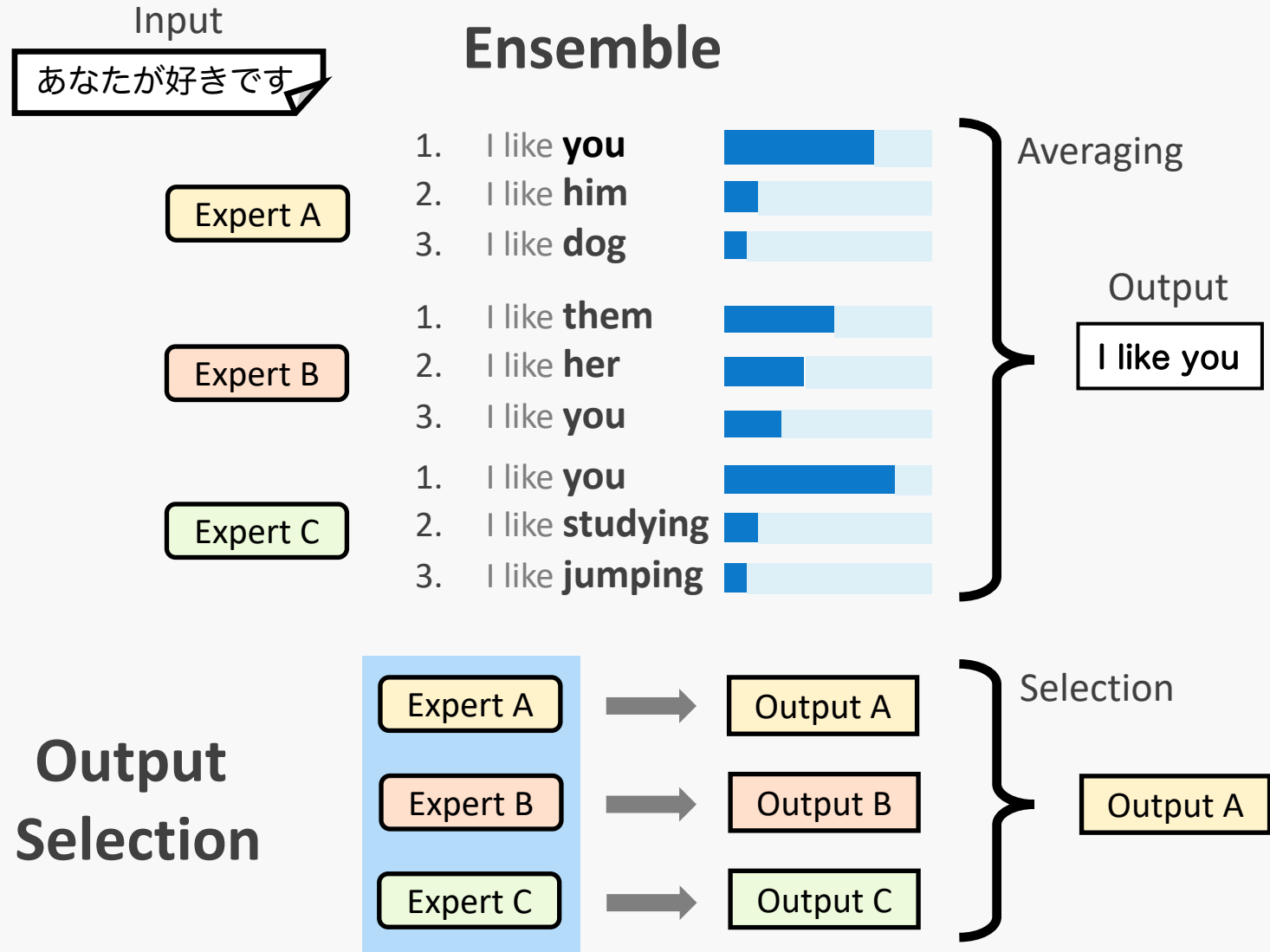    - Investigated effective collaboration methods in DEMs

MDM

DEMs

Mixed Data → Multi-domain Model

Domain A → Expert A
Domain B → Expert B
Domain C → Expert C

# How multiple models collaborate in DEMs

- Ensemble

- Output selection
  - Quality Estimation (QE)
  - Minimum Bayes Risk (MBR)

Input

あなたが好きです
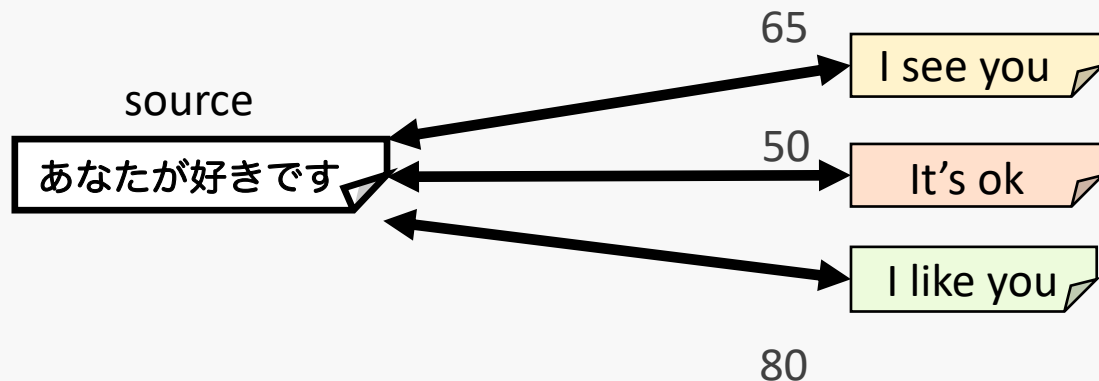
**Ensemble**

Expert A

1. I like **you**
2. I like **him**
3. I like **dog**

Averaging

Expert B

1. I like **them**
2. I like **her**
3. I like **you**

Output

I like you

Expert C

1. I like **you**
2. I like **studying**
3. I like **jumping**

**Output Selection**

| Expert A | → | Output A |
| Expert B | → | Output B |
| Expert C | → | Output C |

Selection

Output A

# Quality Estimation (QE)

$$\arg \max_{C_i \in C} \text{QE} \left( Src, C_i \right)$$

# Quality Estimation (QE)
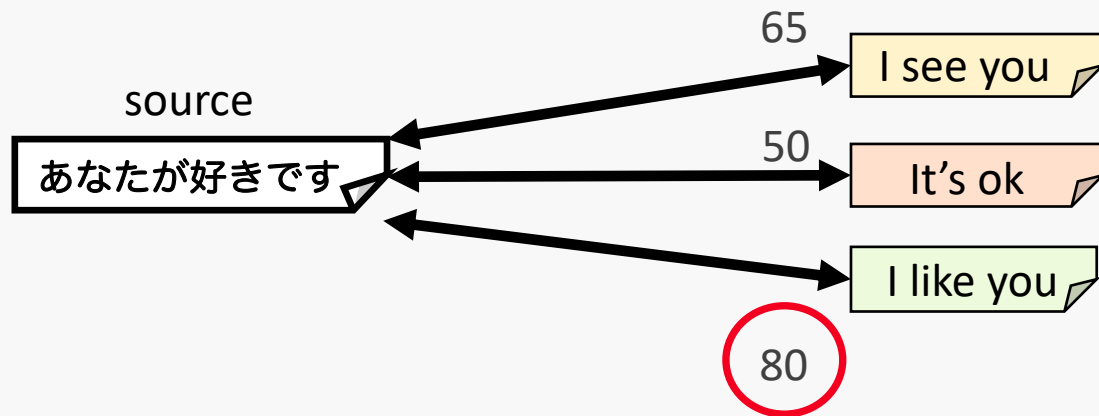
$$\arg\max_{C_i \in C} \boxed{\text{QE}\ (Src,\ C_i)}$$

Calculate quality estimation metric score
we used MS-COMET-QE-22 [Kocmi+ '22]

source

あなたが好きです

65 → I see you

50 → It's ok

I like you
80

# Quality Estimation (QE)

$$\underset{C_i \in C}{\arg\max} \text{QE} (Src, C_i)$$

Select the one with the highest score

65
I see you

source
あなたが好きです

50
It's ok

I like you

80

# Quality Estimation (QE)

$$\underset{C_i \in C}{\arg\max} \, \mathrm{QE} \, (Src, \, C_i)$$

Select the candidate translation
with the highest
quality estimation score

65

I see you

source

あなたが好きです

50

It's ok

I like you

80

What you need
- Source
- Each expert model's outputs

➡ Practical method

# Minimum Bayes Risk (MBR)

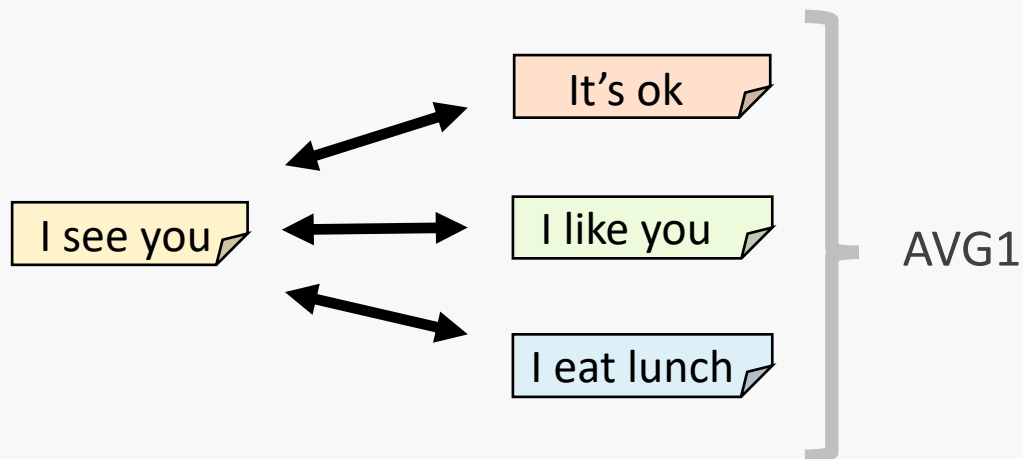$$\arg\max_{C_i \in C} \frac{1}{|C|} \sum_{j \neq i} \text{EVAL}\,(Src,\, C_i,\, C_j)$$

# Minimum Bayes Risk (MBR)

$$\arg\max_{C_i \in C} \boxed{\frac{1}{|C|} \sum_{j \neq i} \text{EVAL}\,(Src,\, C_i,\, \boxed{C_j})}$$

pseudo-references

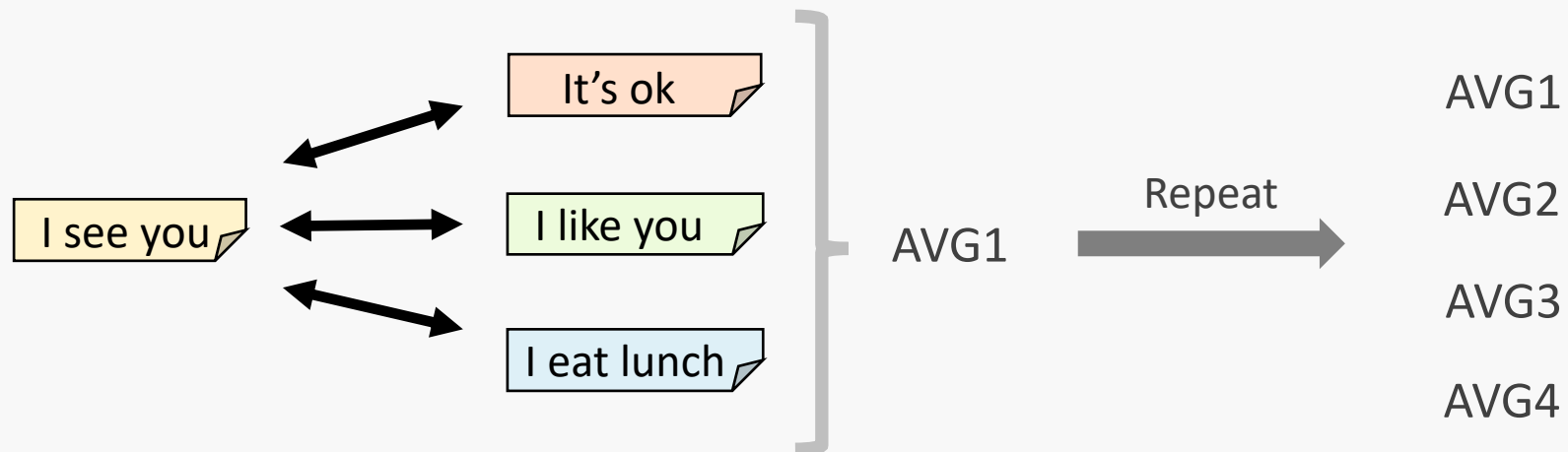Evaluate against all other candidate translations
and take an average score.
We used MS-COMET-22 [Kocmi+ '22]

# Minimum Bayes Risk (MBR)

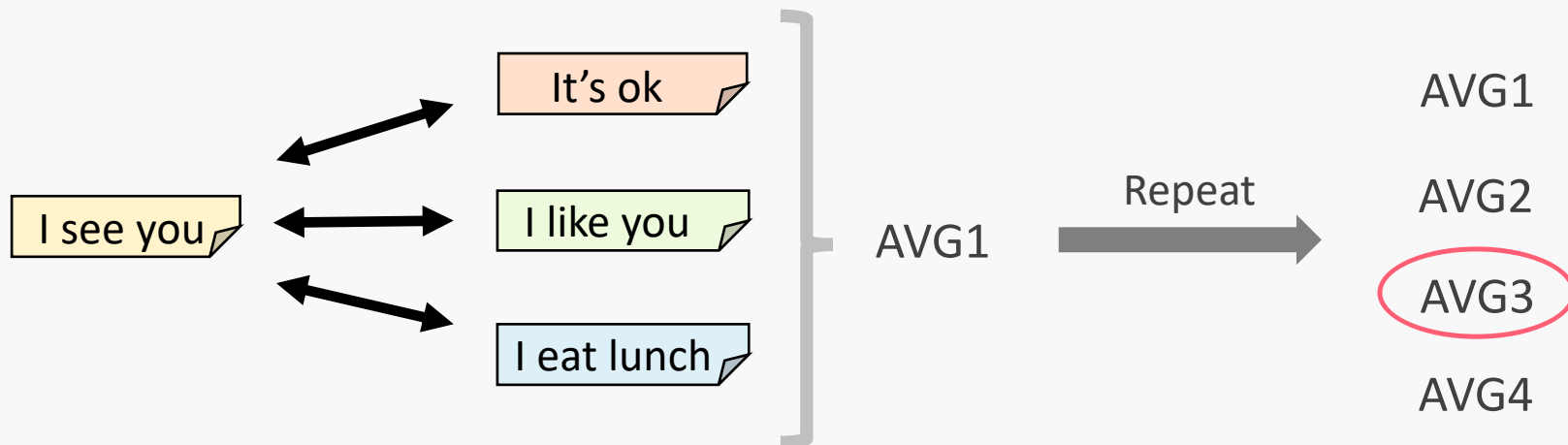$$\arg\max_{C_i \in C} \frac{1}{|C|} \sum_{j \neq i} \text{EVAL} \left( Src, C_i, C_j \right)$$

Repeat the AVG score calculation for all candidates

# Minimum Bayes Risk (MBR)

$$\arg\max_{C_i \in C} \frac{1}{|C|} \sum_{j \neq i} \text{EVAL} \left( Src, C_i, C_j \right)$$

Select the one with the highest score

# Minimum Bayes Risk (MBR)

$$\arg\max_{C_i \in C} \frac{1}{|C|} \sum_{j \neq i} \text{EVAL}\ (Src,\ C_i,\ C_j)$$
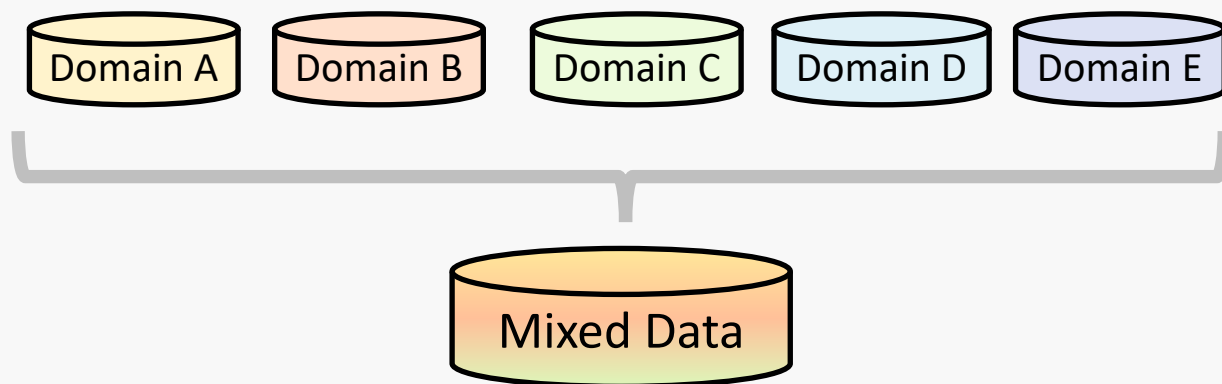
↓

Select a consensus output

What you need
- Source
- Each expert model's outputs

➡ Practical method

# Experimental settings

- Model: Transformer (90M, 290M, 1B)

- Dataset (En-Ja and Ja-En)
  - Pre-train: JParaCrawl v3.0
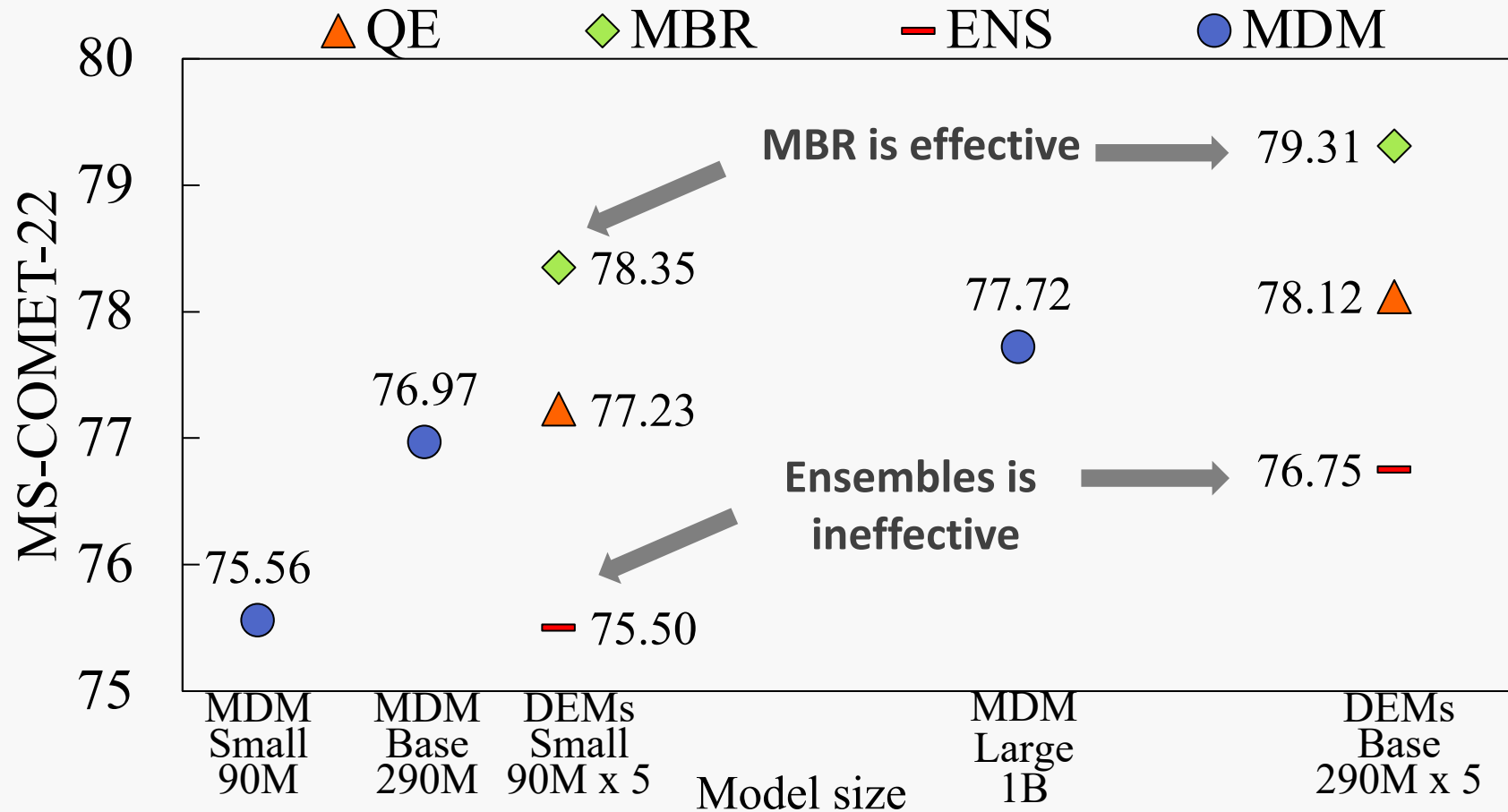  - Fine-tuning: five-specific domain

| Name | Params | Encoder & Decoder | | | |
|---|---|---|---|---|---|
| | | layers | $d_{model}$ | $d_{ffn}$ | heads |
| SMALL | 90M | 6 | 512 | 2048 | 8 |
| BASE | 290M | 6 | 1024 | 4096 | 16 |
| LARGE | 1B | 6 | 2048 | 8192 | 32 |

Model configurations

| Dataset | #Sent Pairs |
|---|---|
| JParaCrawl v3.0 | 25.7M |
| The Kyoto Free Translation Task (KFTT) | 440k |
| Japanese-English Legal Parallel Corpus (LAW) | 260k |
| TED talks (TED) | 225k |
| Asian Scientific Paper Excerpt Corpus (ASPEC) | 200k |
| The Business Scene Dialogue corpus (BSD) | 20k |

Data information

Domain A  Domain B  Domain C  Domain D  Domain E

Mixed Data

# Result: output selection from small experts is effective



△ QE    ◇ MBR    — ENS    ● MDM

MS-COMET-22

MBR is effective → 79.31

78.35

77.72    78.12

76.97

77.23

Ensembles is ineffective → 76.75

75.56

75.50

| MDM Small 90M | MDM Base 290M | DEMs Small 90M x 5 | MDM Large 1B | DEMs Base 290M x 5 |

Model size

※ Ja-En setting (same trend in En-Ja)

# Summary: DEMs can be a hopeful direction in multi-domain

- The performance of 90M x 5 models was comparable to the 1B model

- Collaboration by MBR is effective, especially MBR
  - ※ The outputs of the small experts must include an output comparable to the output of the large model
  - ※ Selection model (e.g., MS-COMET-22) should address multi-domain