

Multi-Stage Multi-Modal Pre-Training For Automatic Speech Recognition

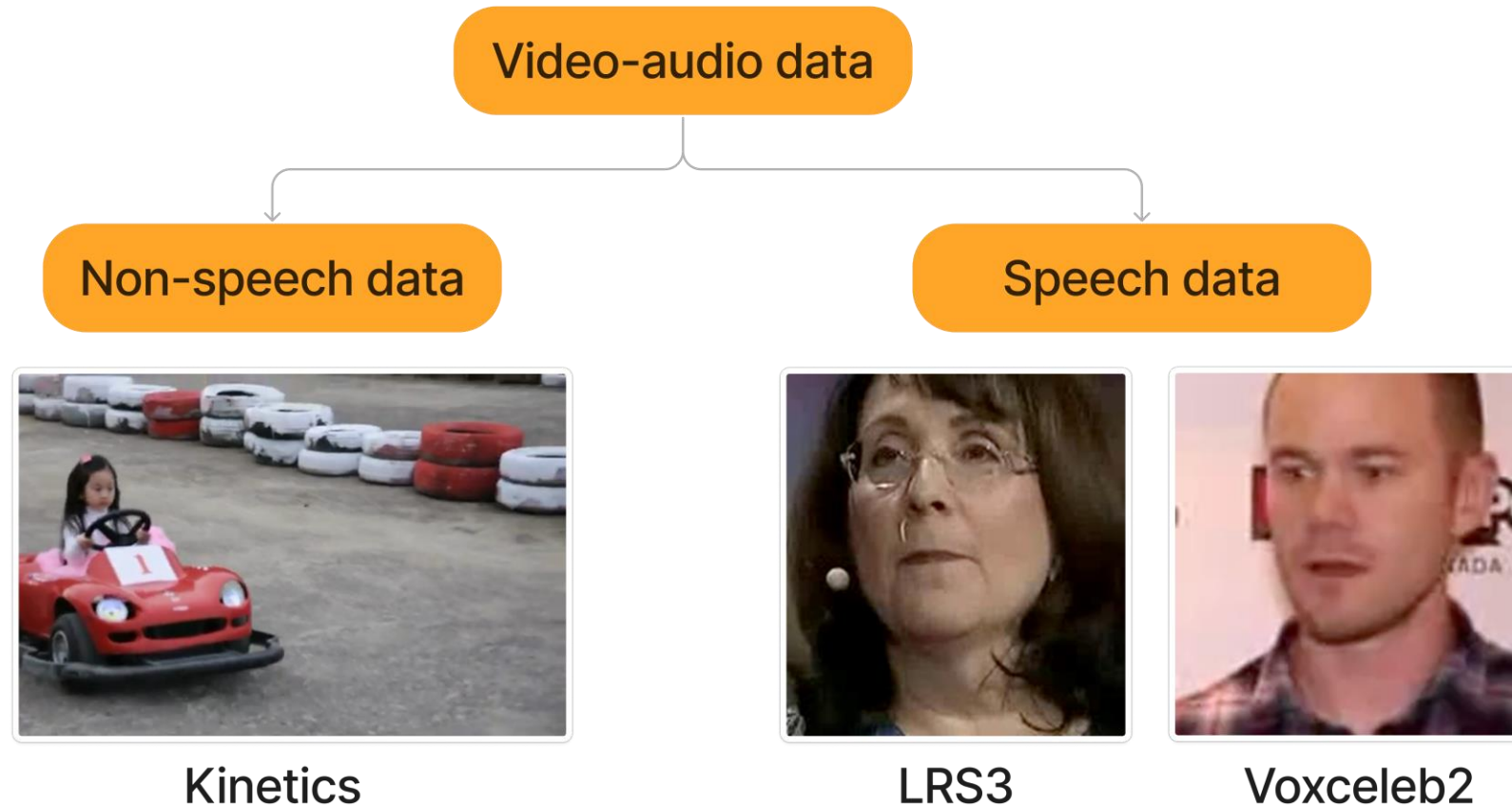
Yash Jain¹, David Chan², Pranav Dheram³, Aparna Khare³, Olabanji Shonibare³, Venkatesh Ravichandran³, Shalini Ghosh³



Multi-Stage Multi-Modal Pre-Training For Automatic Speech Recognition

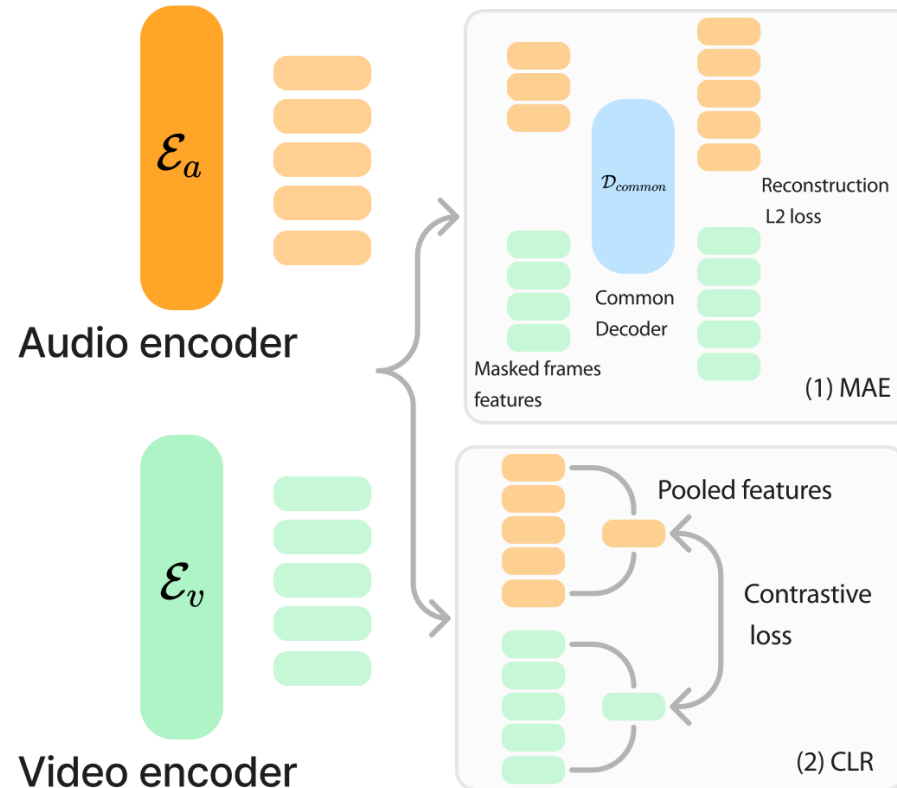
- Large scale evaluation of audio-visual pre-training methods with varying characteristics.
- Evaluate them on Speech recognition task and SUPERB benchmark.
- Introduce a novel mid-training stage between pre-training and fine-tuning for better feature alignment.

Multi-Stage Multi-Modal Pre-Training For Automatic Speech Recognition



Multi-Stage Multi-Modal Pre-Training For Automatic Speech Recognition

1. Masked Autoencoder (MAE)
2. Contrastive Learning (CLR)
3. MAE + CLR



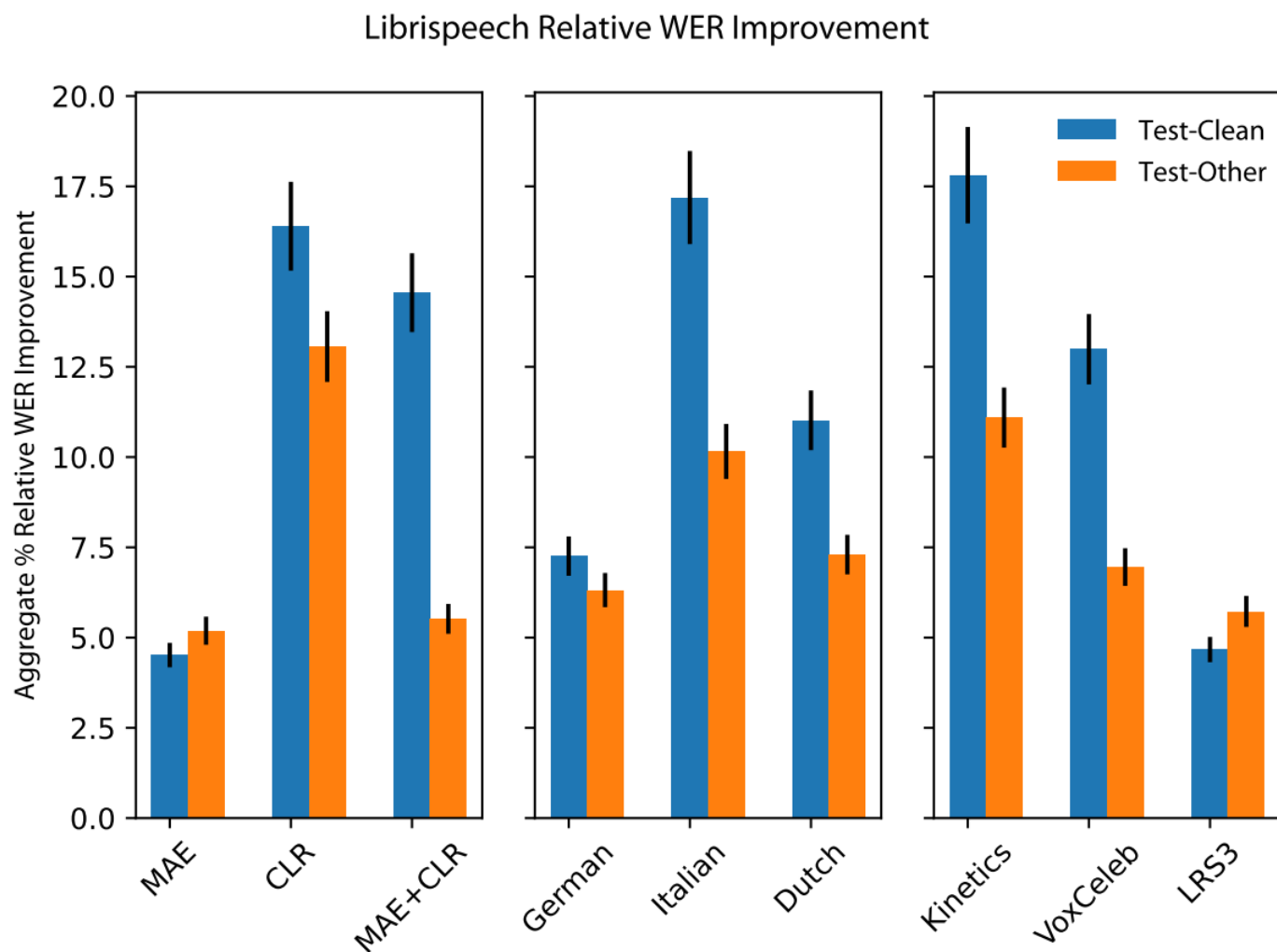
Multi-Stage Multi-Modal Pre-Training For Automatic Speech Recognition

- Automatic Speech Recognition
 - Librispeech dataset
- SUPERB benchmark
 - Keyword spotting
 - Intent Classification
 - Phoneme Recognition
 - Speaker Diarization

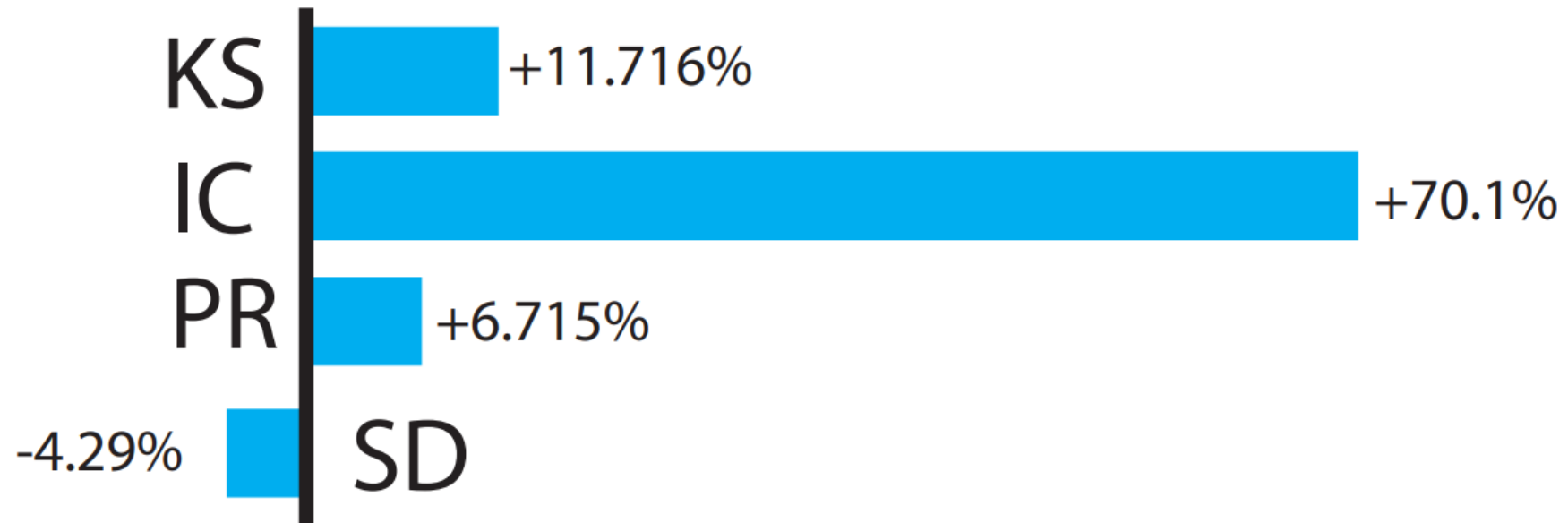
Multi-Stage Multi-Modal Pre-Training For Automatic Speech Recognition



Results: Automatic Speech Recognition



Results: SUPERB benchmark



More Results in the paper!

| Method | PT | MT | en-de ↓ | | en-it ↓ | | en-nl ↓ | |
|-----------------|------|----|-------------|--------------|-------------|--------------|-------------|--------------|
| | | | Test-clean | Test-other | Test-clean | Test-other | Test-clean | Test-other |
| No Pre-training | None | - | 6.84 ± 0.22 | 12.91 ± 0.47 | 6.84 | 12.91 | 6.84 | 12.91 |
| MAE | K600 | - | 7.54 | 13.88 | 7.54 | 13.88 | 7.54 | 13.88 |
| | K600 | ✓ | <u>5.69</u> | <u>11.34</u> | <u>5.95</u> | <u>12.55</u> | <u>5.73</u> | <u>12.04</u> |
| | VC2 | - | 5.28 | 11.51 | 5.28 | 11.51 | 5.28 | 11.51 |
| | VC2 | ✓ | <u>5.11</u> | <u>11.12</u> | <u>5.56</u> | <u>10.42</u> | 5.64 | 12.46 |
| | LRS3 | - | 4.73 | 10.27 | 4.73 | 10.27 | 4.73 | 10.27 |
| | LRS3 | ✓ | 5.61 | 10.85 | 4.21 | 9.53 | 5.32 | 10.33 |
| CLR | K600 | - | 6.85 | 12.92 | 6.85 | 12.92 | 6.85 | 12.92 |
| | K600 | ✓ | <u>5.02</u> | <u>10.85</u> | <u>4.72</u> | <u>10.62</u> | <u>4.65</u> | <u>10.41</u> |
| | VC2 | - | 6.47 | 12.42 | 6.47 | 12.42 | 6.47 | 12.42 |
| | VC2 | ✓ | <u>6.43</u> | <u>12.31</u> | <u>5.1</u> | <u>10.61</u> | <u>4.62</u> | <u>10.77</u> |
| | LRS3 | - | 6.35 | 12.12 | 6.35 | 12.12 | 6.35 | 12.12 |
| | LRS3 | ✓ | 6.74 | <u>10.59</u> | <u>5.84</u> | <u>11.33</u> | <u>6.01</u> | <u>10.13</u> |
| MAE + CLR | K600 | - | 5.56 | 11.91 | 5.56 | 11.91 | 5.56 | 11.91 |
| | K600 | ✓ | <u>5.02</u> | <u>11.68</u> | <u>5.23</u> | <u>11.37</u> | 6.39 | 12.03 |
| | VC2 | - | 6.75 | 12.11 | 6.75 | 12.11 | 6.75 | 12.11 |
| | VC2 | ✓ | <u>5.36</u> | <u>11.22</u> | <u>4.77</u> | <u>10.84</u> | <u>5.03</u> | <u>10.73</u> |
| | LRS3 | - | 7.51 | 12.54 | 7.51 | 12.54 | 7.51 | 12.54 |
| | LRS3 | ✓ | <u>7.16</u> | <u>12.29</u> | <u>5.08</u> | <u>11.13</u> | <u>6.17</u> | <u>12.32</u> |

Thank you!

Meet us at the poster!