



國立臺灣大學  
National Taiwan University

# Speech Representation Learning through Self-supervised Pretraining and Multi-task Finetuning

Yi-Chen Chen, Shu-wen Yang, Cheng-Kuang Lee, Simon See, Hung-yi Lee

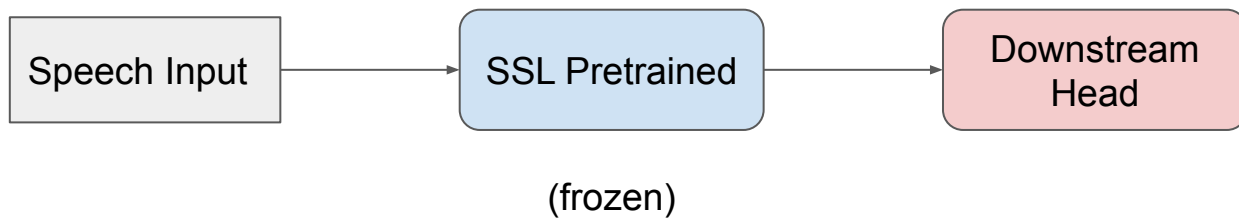


# Motivation - Representation Learning via Pretraining

- For example,
  - ImageNet pretraining in CV -> object detection
  - BERT in NLP -> question answering
  - wav2vec in speech -> automatic speech recognition

# Motivation - Representation Learning via SSL Pretraining

- Self-supervised Learning (SSL)
  - Generative losses
    - APC, Mockingjay, Tera, DeCoAR, ...
  - Discriminative losses
    - CPC, wav2vec (2.0), HuBERT, ...
  - Multiple losses
    - PASE(+), ...



# Motivation - Representation Learning via Multi-task Learning

- General representation learning for a variety of speech processing tasks
- Supervised multi-task learning (MTL) is to train a shared model on various downstream tasks.
- There has not been a systematic study of general representation learning models trained by MTL of various speech processing tasks.
- We want to investigate if MTL on various downstream tasks can further improve the representations from SSL.
- SUPERB Benchmark (Unconstrained)

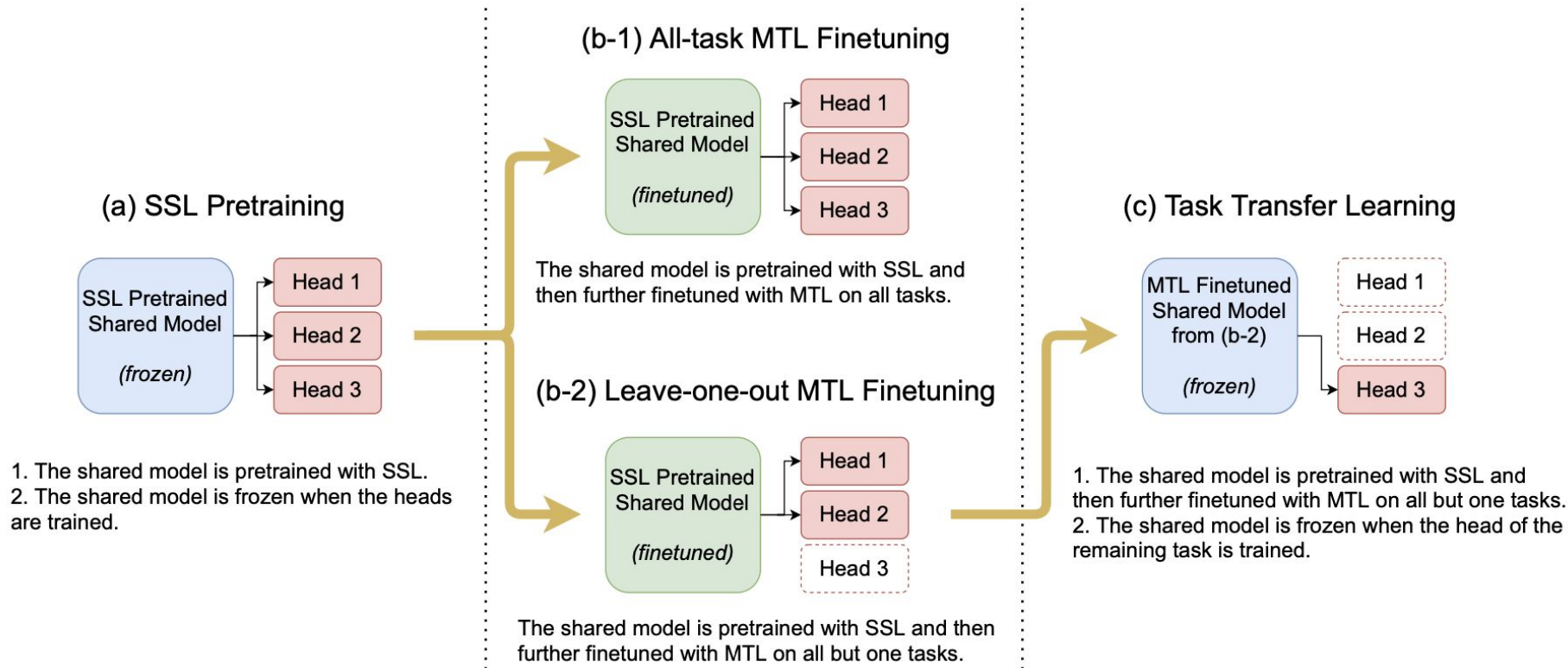
# SUPERB

- Since no downstream model training is required in QbE, we only perform MTL experiments and compare the results on the other nine tasks.
  - Content
    - Phoneme Recognition (PR)
    - Automatic Speech Recognition (ASR)
    - Keyword Spotting (KS)
  - Speaker
    - Speaker Identification (SID)
    - Automatic Speaker Verification (ASV)
    - Speaker Diarization (SD)
  - Semantics
    - Intent Classification (IC)
    - Slot Filling (SF)
  - Paralinguistics
    - Emotion Recognition (ER)

# Experimental Setup

- The SSL pretraining approach used in experiments:
  - HuBERT achieves the overall best performance on SUPERB.
  - We use a weighted sum of hidden representations of all layers in the HuBERT model as the representations for downstream heads, as in SUPERB.
- The model architecture and implementation details:
  - HuBERT Base is adopted as our shared model architecture.
  - For task-specific head architectures, we simply follow the settings in SUPERB.
  - The code is released: <https://github.com/s3prl/s3prl/tree/multi-task-distributed>

# Training Scenarios



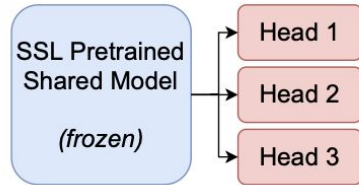
# Experimental Results

Scenario	Tasks for MTL Finetuning	ASR	PR	SF		SD	ER	IC	KS	ASV	SID
		WER↓	PER↓	F1↑	CER↓	DER↓	ACC↑	ACC↑	ACC↑	EER↓	ACC↑
(a) SSL	N/A	6.42	5.41	88.53	25.20	5.88	64.24	98.34	96.30	5.11	81.42
(b-1) SSL+MTL	all	6.22	3.61	87.56	26.76	4.93	67.28	99.60	97.34	6.76	90.86
(b-2): SSL+MTL	all but ASR	X	<u>3.63</u>	<u>87.28</u>	<u>27.11</u>	<b>4.89</b>	<u>65.07</u>	<b>99.63</b>	<b>97.57</b>	<u>7.78</u>	<u>90.69</u>
	all but PR	<u>6.79</u>	X	<u>86.94</u>	<u>27.66</u>	<b>4.81</b>	<u>66.73</u>	<b>99.66</b>	<b>97.44</b>	<u>7.94</u>	<b>91.16</b>
	all but SF	<b>6.10</b>	<b>3.39</b>	X	X	<b>4.73</b>	<u>65.71</u>	<u>99.58</u>	<u>97.18</u>	<u>7.61</u>	<u>90.70</u>
	all but SD	<u>6.28</u>	<b>3.54</b>	<b>87.94</b>	<b>26.31</b>	X	<u>66.73</u>	<b>99.63</b>	<u>97.11</u>	<u>7.49</u>	<u>90.79</u>
	all but ER	<b>6.17</b>	<b>3.40</b>	<u>87.45</u>	<u>26.90</u>	<b>4.77</b>	X	<u>99.55</u>	<u>97.27</u>	<u>7.19</u>	<u>90.51</u>
	all but IC	<b>6.13</b>	<b>3.34</b>	<b>87.65</b>	<u>26.94</u>	<b>4.78</b>	<u>66.08</u>	X	<u>97.27</u>	<b>6.74</b>	<u>90.55</u>
	all but KS	<b>6.17</b>	<b>3.55</b>	<b>87.83</b>	<u>26.88</u>	<b>4.91</b>	<u>66.27</u>	<b>99.71</b>	X	<u>7.86</u>	<u>90.67</u>
	all but ASV	<b>5.90</b>	<b>2.79</b>	<b>87.88</b>	<b>26.52</b>	<b>3.61</b>	<u>64.88</u>	<u>99.58</u>	<b>97.44</b>	X	<u>85.06</u>
	all but SID	<b>5.95</b>	<b>3.25</b>	<u>87.33</u>	<u>27.39</u>	<b>4.50</b>	<b>68.66</b>	<u>99.55</u>	<u>97.27</u>	<u>9.00</u>	X
(c) Task Transfer	N/A	6.27	5.79	88.14	26.24	5.80	64.24	97.42	96.33	7.55	62.05



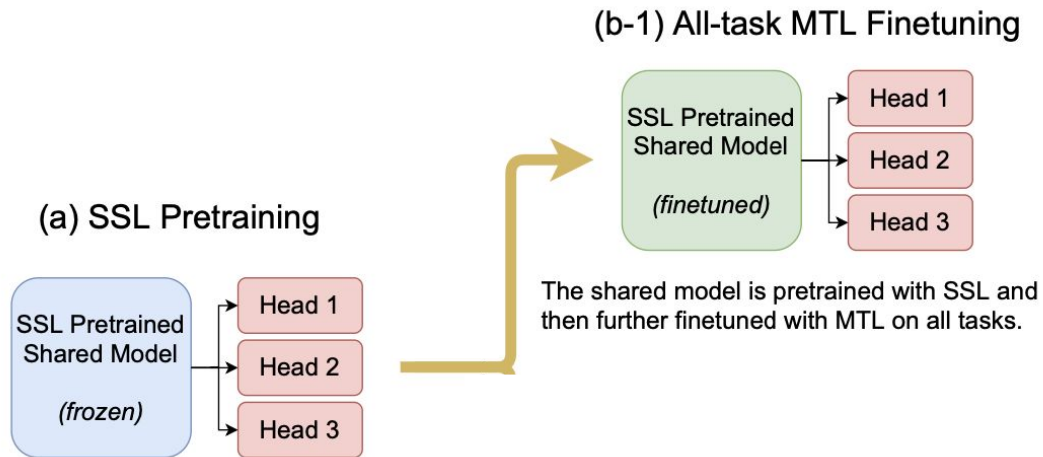
# MTL Scenarios - SSL Pretraining

(a) SSL Pretraining



1. The shared model is pretrained with SSL.
2. The shared model is frozen when the heads are trained.

# MTL Scenarios - SSL Pretraining



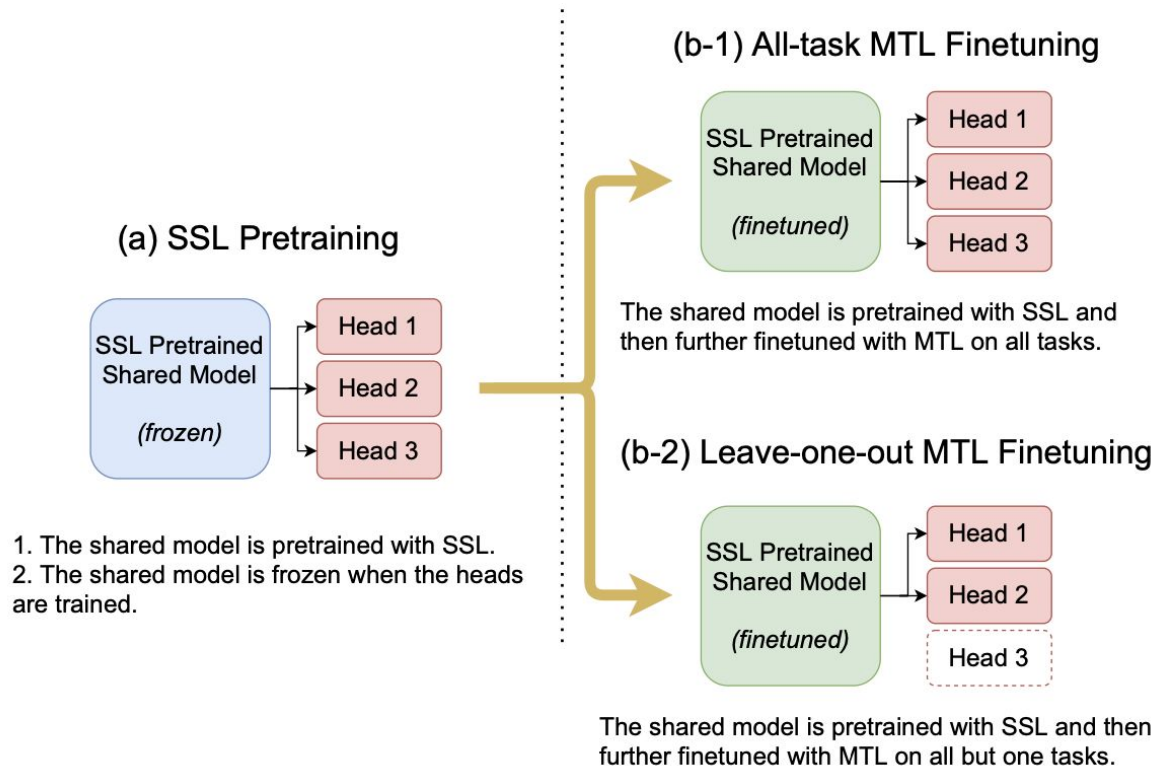
1. The shared model is pretrained with SSL.
2. The shared model is frozen when the heads are trained.

# Experimental Results - All-task MTL Finetuning

Scenario	Tasks for MTL Finetuning	<b>ASR</b>	<b>PR</b>	<b>SF</b>		<b>SD</b>	<b>ER</b>	<b>IC</b>	<b>KS</b>	<b>ASV</b>	<b>SID</b>
		<b>WER↓</b>	<b>PER↓</b>	<b>F1↑</b>	<b>CER↓</b>	<b>DER↓</b>	<b>ACC↑</b>	<b>ACC↑</b>	<b>ACC↑</b>	<b>EER↓</b>	<b>ACC↑</b>
(a) SSL	N/A	6.42	5.41	88.53	25.20	5.88	64.24	98.34	96.30	5.11	81.42
(b-1) SSL+MTL	all	6.22	3.61	87.56	26.76	4.93	67.28	99.60	97.34	6.76	90.86

- Serve as a strong baseline for SSL pretraining or other representation learning approaches.
- How to select the model checkpoint based on validation scores of tasks?
  - ASV suffers from overfitting.

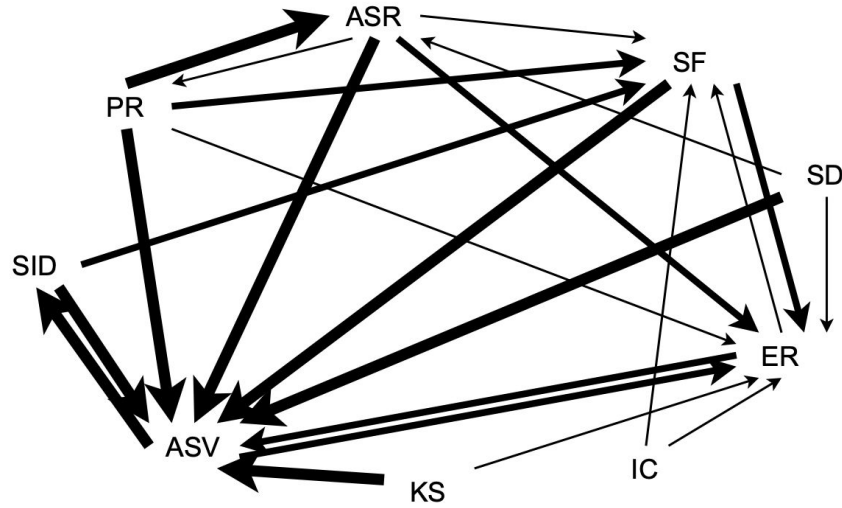
# MTL Scenarios - SSL Pretraining



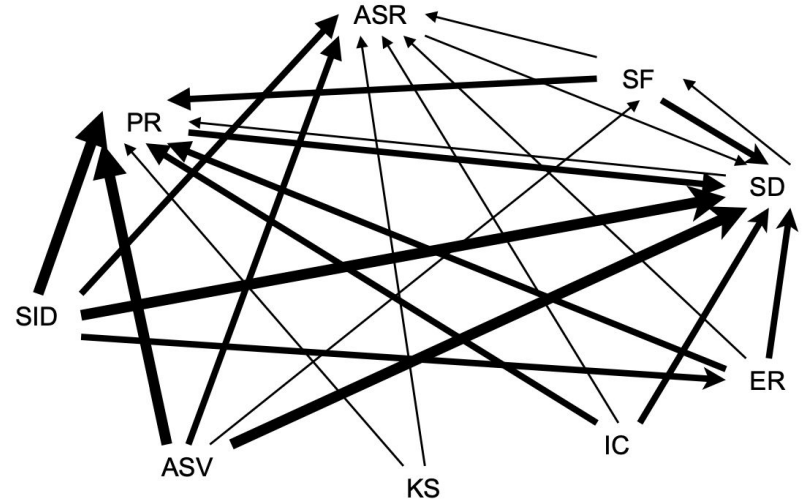
# Experimental Results - Removing One Task In MTL

Scenario	Tasks for MTL Finetuning	ASR	PR	SF		SD	ER	IC	KS	ASV	SID
		WER↓	PER↓	F1↑	CER↓	DER↓	ACC↑	ACC↑	ACC↑	EER↓	ACC↑
(b-1) SSL+MTL	all	6.22	3.61	87.56	26.76	4.93	67.28	99.60	97.34	6.76	90.86
(b-2): SSL+MTL	all but ASR	X	<u>3.63</u>	<u>87.28</u>	<u>27.11</u>	<b>4.89</b>	<u>65.07</u>	<b>99.63</b>	<b>97.57</b>	<u>7.78</u>	<u>90.69</u>
	all but PR	<u>6.79</u>	X	<u>86.94</u>	<u>27.66</u>	<b>4.81</b>	<u>66.73</u>	<b>99.66</b>	<b>97.44</b>	<u>7.94</u>	<b>91.16</b>
	all but SF	<b>6.10</b>	<b>3.39</b>	X	X	<b>4.73</b>	<u>65.71</u>	<u>99.58</u>	<u>97.18</u>	<u>7.61</u>	<u>90.70</u>
	all but SD	<u>6.28</u>	<b>3.54</b>	<b>87.94</b>	<b>26.31</b>	X	<u>66.73</u>	<b>99.63</b>	<u>97.11</u>	<u>7.49</u>	<u>90.79</u>
	all but ER	<b>6.17</b>	<b>3.40</b>	<u>87.45</u>	<u>26.90</u>	<b>4.77</b>	X	<u>99.55</u>	<u>97.27</u>	<u>7.19</u>	<u>90.51</u>
	all but IC	<b>6.13</b>	<b>3.34</b>	<b>87.65</b>	<u>26.94</u>	<b>4.78</b>	<u>66.08</u>	X	<u>97.27</u>	<b>6.74</b>	<u>90.55</u>
	all but KS	<b>6.17</b>	<b>3.55</b>	<b>87.83</b>	<u>26.88</u>	<b>4.91</b>	<u>66.27</u>	<b>99.71</b>	X	<u>7.86</u>	<u>90.67</u>
	all but ASV	<b>5.90</b>	<b>2.79</b>	<b>87.88</b>	<b>26.52</b>	<b>3.61</b>	<u>64.88</u>	<u>99.58</u>	<b>97.44</b>	X	<u>85.06</u>
	all but SID	<b>5.95</b>	<b>3.25</b>	<u>87.33</u>	<u>27.39</u>	<b>4.50</b>	<b>68.66</b>	<u>99.55</u>	<u>97.27</u>	<u>9.00</u>	X

# Experimental Results - Removing One Task In MTL



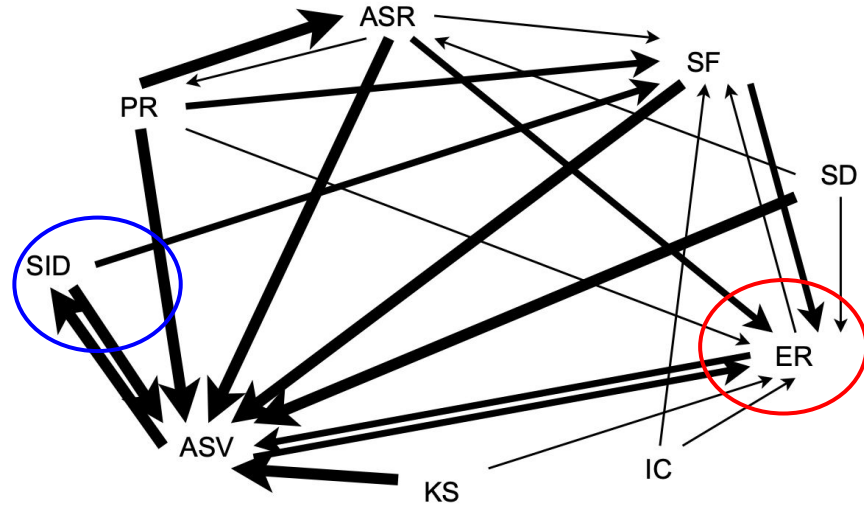
The improvement relations of tasks.



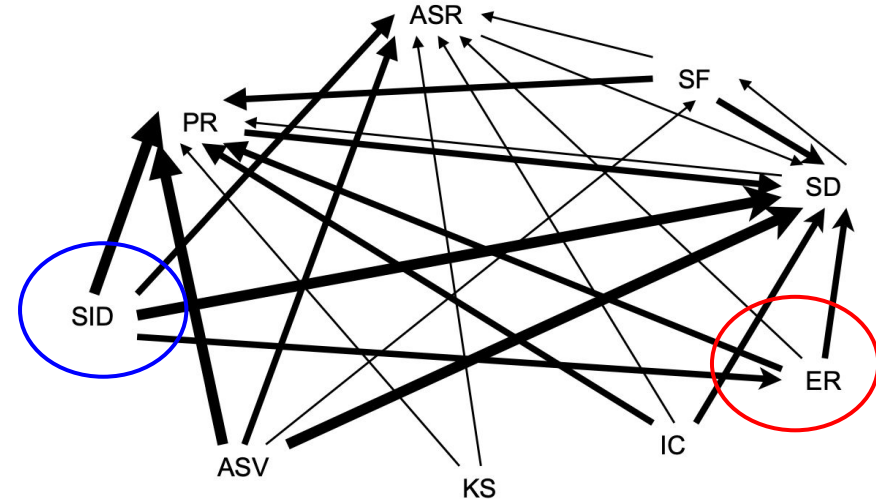
The hurt relations of tasks.

- If we focus on a certain primary task, we may select proper auxiliary tasks to assist the primary task based on these MTL experimental results.

# Experimental Results - Removing One Task In MTL



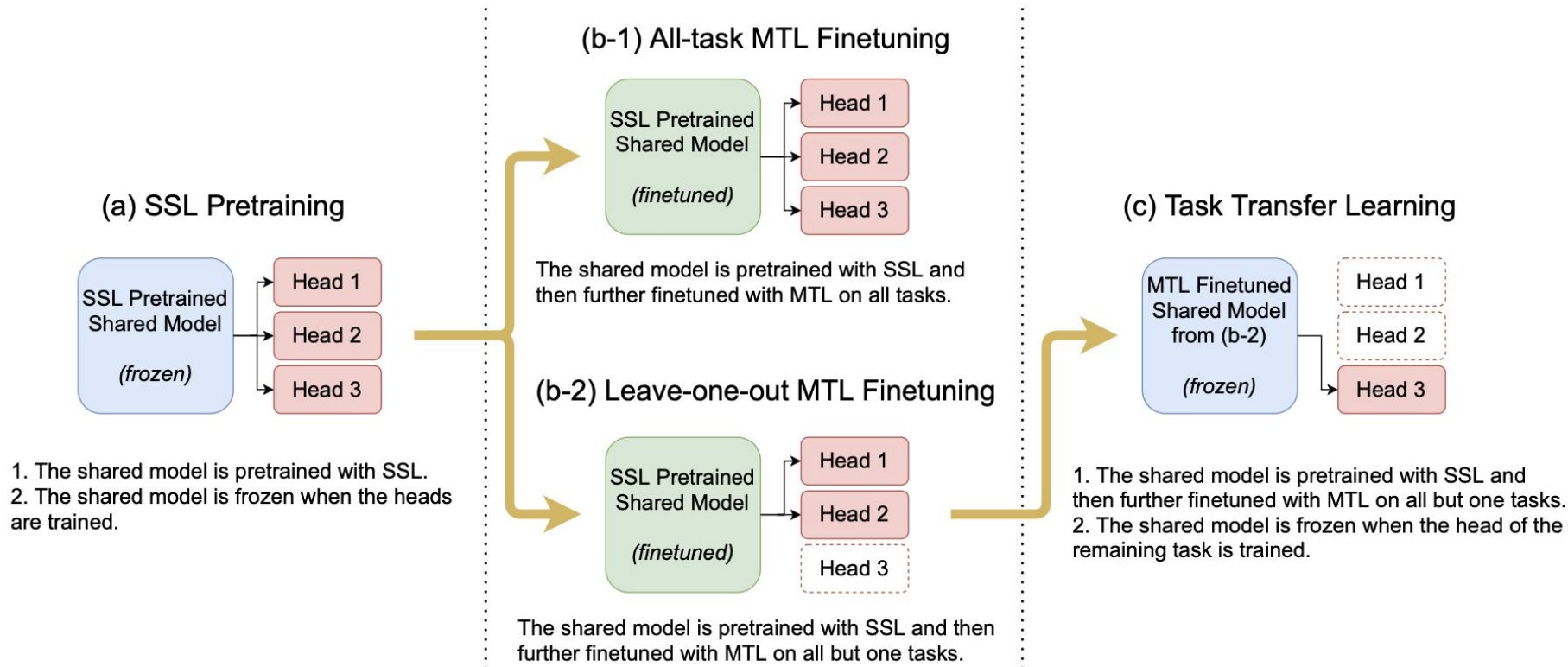
The improvement relations of tasks.



The hurt relations of tasks.

- If we focus on a certain primary task, we may select proper auxiliary tasks to assist the primary task based on these MTL experimental results.

# MTL Scenarios - SSL Pretraining





# Experimental Results - Task Transfer Learning

Scenario	Tasks for MTL Finetuning	<b>ASR</b>	<b>PR</b>	<b>SF</b>		<b>SD</b>	<b>ER</b>	<b>IC</b>	<b>KS</b>	<b>ASV</b>	<b>SID</b>
		<b>WER↓</b>	<b>PER↓</b>	<b>F1↑</b>	<b>CER↓</b>	<b>DER↓</b>	<b>ACC↑</b>	<b>ACC↑</b>	<b>ACC↑</b>	<b>EER↓</b>	<b>ACC↑</b>
(a) SSL	N/A	6.42	5.41	88.53	25.20	5.88	64.24	98.34	96.30	5.11	81.42
(b-1) SSL+MTL	all	6.22	3.61	87.56	26.76	4.93	67.28	99.60	97.34	6.76	90.86

(c) Task Transfer	N/A	6.27	5.79	88.14	26.24	5.80	64.24	97.42	96.33	7.55	62.05
-------------------	-----	------	------	-------	-------	------	-------	-------	-------	------	-------

- Whether a task is involved in MTL is crucial to the performance of this task.

# Conclusion and Discussion

- In this work, we investigate different training scenarios of supervised MTL as a speech representation learning approach along with SSL pretraining on a benchmark with various speech processing tasks.
  - We analyze the generalizability of representations learned with supervised MTL empirically.
- The performance of MTL is dependent on many factors.
  - the amount of data
  - task relationships
  - noise
  - These factors should be isolated and investigated with more analyses.
- Future directions:
  - exploring a better method to select the model checkpoint with MTL
  - more in-depth research of MTL and its optimization on speech processing tasks
  - trying to train the shared model with both SSL, MTL and self training simultaneously as a semi-supervised representation learning approach