

# User behavior modeling

With Multi-Modal Attention-based  
Hierarchical Graph Neural Network  
(MM-AHGNN)

## References:

1. *A Survey on Knowledge Graphs: Representation, Acquisition and Applications, IEEE Transactions on Neural Networks and Learning Systems 2021*
2. *Object Interaction Recommendation with Multi-Modal Attention-based Hierarchical Graph Neural Network, IEEE International Conference on Big Data 2021*

## Objective:

### コンテンツ:

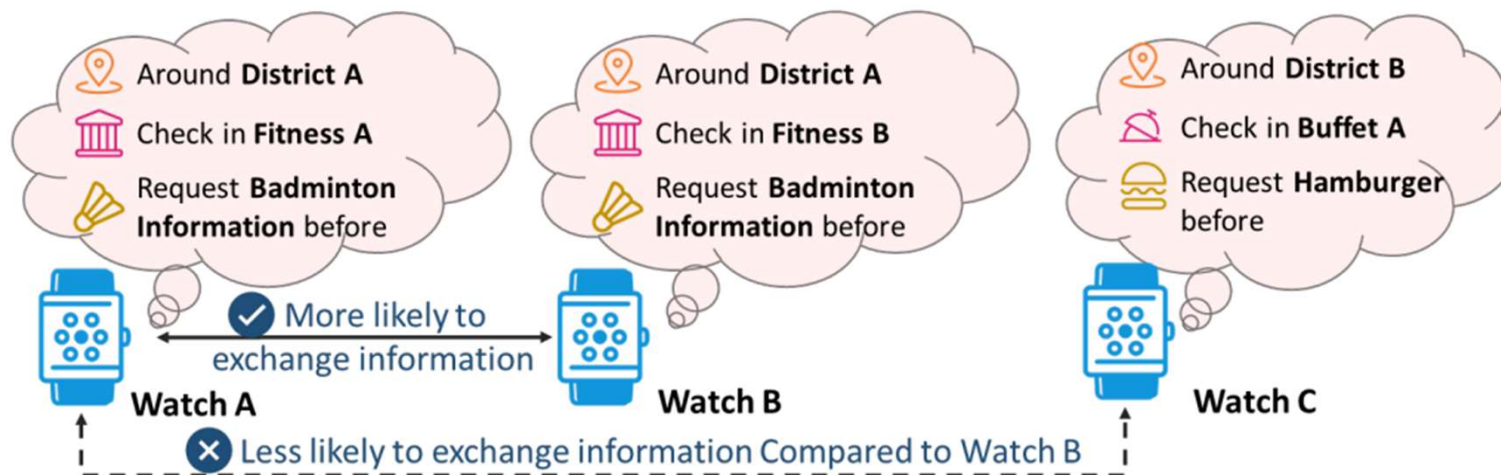
1. 領域紹介及び背景知識.
2. ビジネス側のメリット.
3. 従来手法及び問題点.
4. MM-AHGNN.

# User behavior modeling

---

1. What is user behavior?
2. Knowledge Graphとは？どのような従来手法がある？
3. MM-AHGNN
  - GNN: 図理論で動態数のデータでのモデリング技術
  - MM-AHGNN: GNN をKnowledge Graph領域に応用
  - 性能分析
4. まとめ

# 行動歴史モデリングとは？



実体間の行動から実体関係を推論する手法

例: スマートウォッチユーザーの行動歴史

- ユーザー: {行動1, 行動2, ...}
- ユーザーA, Bの行動パターンが近いよって関連性が高い
- ユーザーA, Cの行動パターンが遠いよって関連性が低い

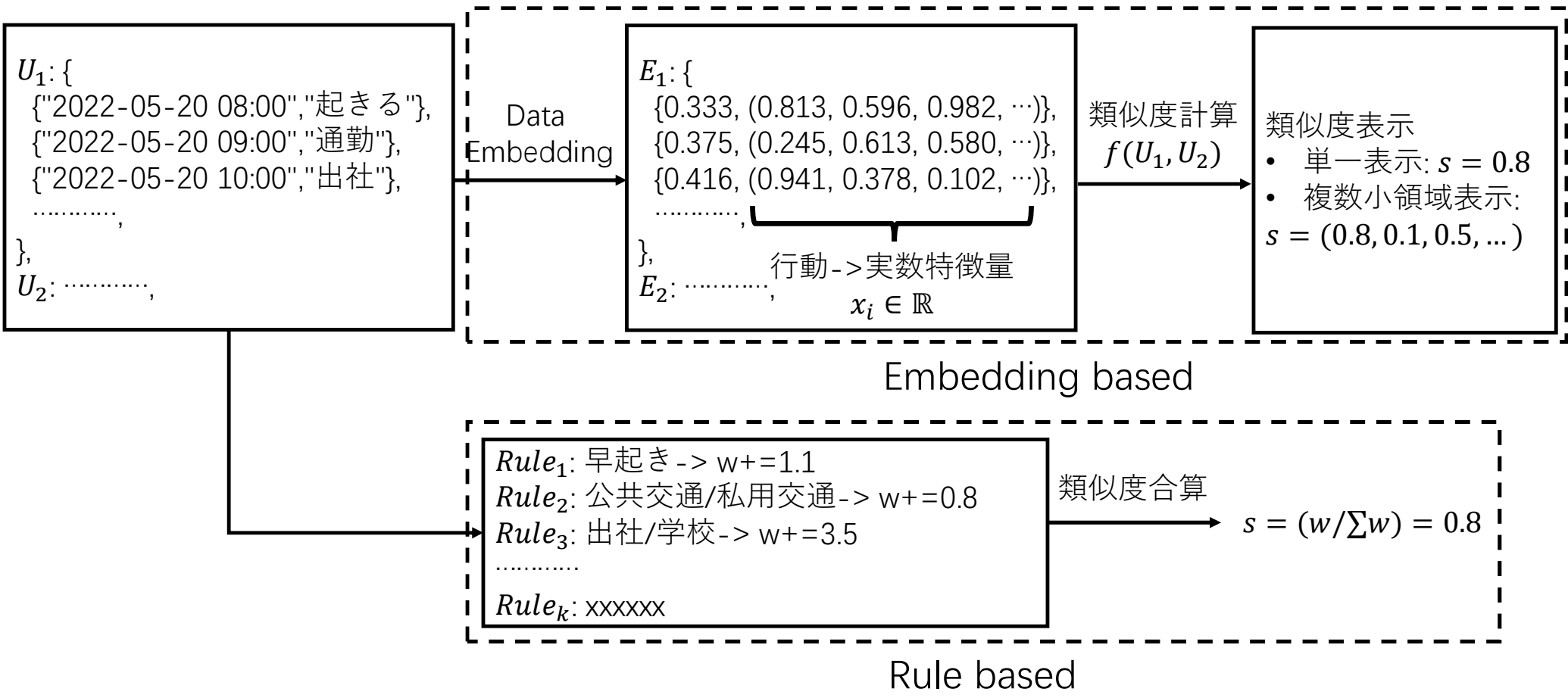
関連性の応用例:

- 興味推論: Aの興味 $\approx$ Bの興味
- 行動予測: A, B {行動1, 行動2}; A{行動3}  $\rightarrow$  B{行動3}
- ユーザーグルーピング: ユーザーを種類別分類できる

# 手順の全体像は？

実体間の行動から実体関係を推論する手法

- 入力:  $U_1\{A_1, \dots, A_n\}, U_2\{A_1, \dots, A_m\}, \dots$
- 出力:  $s_{iht} \in S, s_{iht} = f_i(E_h, E_t)$

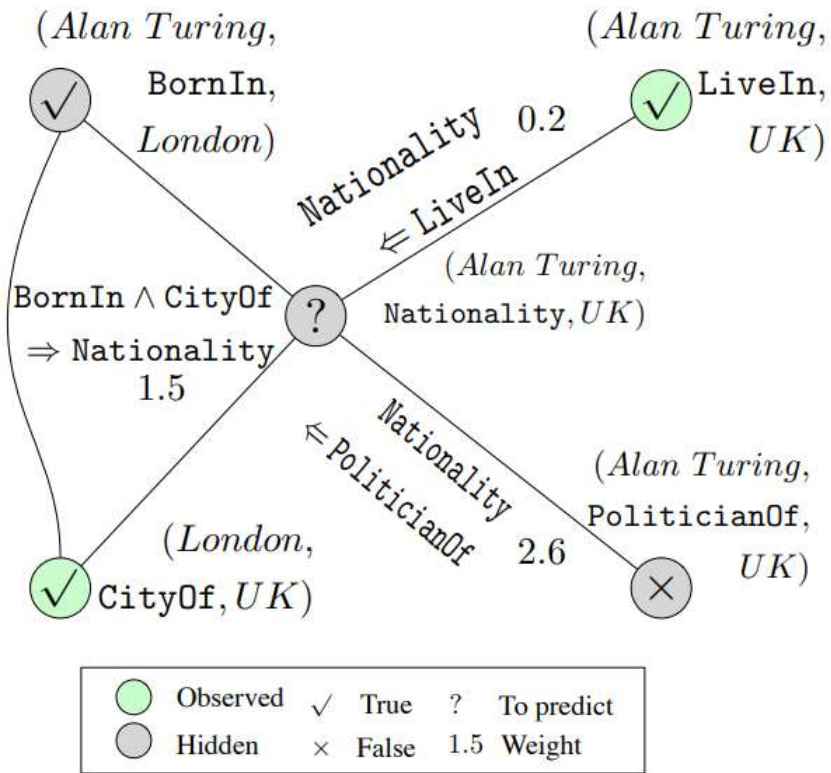


# ユーザー行動歴史モデリング

---

1. 行動歴史モデリングとは？手順の全体像は？
2. Knowledge Graphとは？どのような従来手法がある？
3. MM-AHGNN
  - GNN: 図理論で動態数のデータでのモデリング技術
  - MM-AHGNN: GNN をKnowledge Graph領域に応用
  - 性能分析
4. まとめ

# Knowledge Graphとは？



関連性の応用例:

- 興味推論: Aの興味 $\approx$ Bの興味
- 行動予測: A,B {行動1,行動2}; A{行動3}  $\rightarrow$  B{行動3}
- ユーザーグルーピング: ユーザーを種類別分類できる

Method	Ent. embedding	Rel. embedding	Scoring function $f_r(h, t)$	Constraints/Regularization
TransE [14]	$\mathbf{h}, \mathbf{t} \in \mathbb{R}^d$	$\mathbf{r} \in \mathbb{R}^d$	$-\ \mathbf{h} + \mathbf{r} - \mathbf{t}\ _{1/2}$	$\ \mathbf{h}\ _2 = 1, \ \mathbf{t}\ _2 = 1$
TransH [15]	$\mathbf{h}, \mathbf{t} \in \mathbb{R}^d$	$\mathbf{r}, \mathbf{w}_r \in \mathbb{R}^d$	$-\ (\mathbf{h} - \mathbf{w}_r^\top \mathbf{h} \mathbf{w}_r) + \mathbf{r} - (\mathbf{t} - \mathbf{w}_r^\top \mathbf{t} \mathbf{w}_r)\ _2^2$	$\ \mathbf{h}\ _2 \leq 1, \ \mathbf{t}\ _2 \leq 1$ $ \mathbf{w}_r^\top \mathbf{r}  / \ \mathbf{r}\ _2 \leq \epsilon, \ \mathbf{w}_r\ _2 = 1$
TransR [16]	$\mathbf{h}, \mathbf{t} \in \mathbb{R}^d$	$\mathbf{r} \in \mathbb{R}^k, \mathbf{M}_r \in \mathbb{R}^{k \times d}$	$-\ \mathbf{M}_r \mathbf{h} + \mathbf{r} - \mathbf{M}_r \mathbf{t}\ _2^2$	$\ \mathbf{h}\ _2 \leq 1, \ \mathbf{t}\ _2 \leq 1, \ \mathbf{r}\ _2 \leq 1$ $\ \mathbf{M}_r \mathbf{h}\ _2 \leq 1, \ \mathbf{M}_r \mathbf{t}\ _2 \leq 1$
TransD [50]	$\mathbf{h}, \mathbf{w}_h \in \mathbb{R}^d$ $\mathbf{t}, \mathbf{w}_t \in \mathbb{R}^d$	$\mathbf{r}, \mathbf{w}_r \in \mathbb{R}^k$	$-\ (\mathbf{w}_r \mathbf{w}_h^\top + \mathbf{I})\mathbf{h} + \mathbf{r} - (\mathbf{w}_r \mathbf{w}_t^\top + \mathbf{I})\mathbf{t}\ _2^2$	$\ \mathbf{h}\ _2 \leq 1, \ \mathbf{t}\ _2 \leq 1, \ \mathbf{r}\ _2 \leq 1$ $\ (\mathbf{w}_r \mathbf{w}_h^\top + \mathbf{I})\mathbf{h}\ _2 \leq 1$ $\ (\mathbf{w}_r \mathbf{w}_t^\top + \mathbf{I})\mathbf{t}\ _2 \leq 1$
TransSparse [51]	$\mathbf{h}, \mathbf{t} \in \mathbb{R}^d$	$\mathbf{r} \in \mathbb{R}^k, \mathbf{M}_r(\theta_r) \in \mathbb{R}^{k \times d}$ $\mathbf{M}_r^1(\theta_r^1), \mathbf{M}_r^2(\theta_r^2) \in \mathbb{R}^{k \times d}$	$-\ \mathbf{M}_r(\theta_r)\mathbf{h} + \mathbf{r} - \mathbf{M}_r(\theta_r)\mathbf{t}\ _{1/2}^2$ $-\ \mathbf{M}_r^1(\theta_r^1)\mathbf{h} + \mathbf{r} - \mathbf{M}_r^2(\theta_r^2)\mathbf{t}\ _{1/2}^2$	$\ \mathbf{h}\ _2 \leq 1, \ \mathbf{t}\ _2 \leq 1, \ \mathbf{r}\ _2 \leq 1$ $\ \mathbf{M}_r(\theta_r)\mathbf{h}\ _2 \leq 1, \ \mathbf{M}_r(\theta_r)\mathbf{t}\ _2 \leq 1$ $\ \mathbf{M}_r^1(\theta_r^1)\mathbf{h}\ _2 \leq 1, \ \mathbf{M}_r^2(\theta_r^2)\mathbf{t}\ _2 \leq 1$
TransM [52]	$\mathbf{h}, \mathbf{t} \in \mathbb{R}^d$	$\mathbf{r} \in \mathbb{R}^d$	$-\theta_r \ \mathbf{h} + \mathbf{r} - \mathbf{t}\ _{1/2}$	$\ \mathbf{h}\ _2 = 1, \ \mathbf{t}\ _2 = 1$
ManifoldE [53]	$\mathbf{h}, \mathbf{t} \in \mathbb{R}^d$	$\mathbf{r} \in \mathbb{R}^d$	$-(\ \mathbf{h} + \mathbf{r} - \mathbf{t}\ _2^2 - \theta_r^2)^2$	$\ \mathbf{h}\ _2 \leq 1, \ \mathbf{t}\ _2 \leq 1, \ \mathbf{r}\ _2 \leq 1$
TransF [54]	$\mathbf{h}, \mathbf{t} \in \mathbb{R}^d$	$\mathbf{r} \in \mathbb{R}^d$	$(\mathbf{h} + \mathbf{r})^\top \mathbf{t} + (\mathbf{t} - \mathbf{r})^\top \mathbf{h}$	$\ \mathbf{h}\ _2 \leq 1, \ \mathbf{t}\ _2 \leq 1, \ \mathbf{r}\ _2 \leq 1$
TransA [55]	$\mathbf{h}, \mathbf{t} \in \mathbb{R}^d$	$\mathbf{r} \in \mathbb{R}^d, \mathbf{M}_r \in \mathbb{R}^{d \times d}$	$-(\ \mathbf{h} + \mathbf{r} - \mathbf{t}\ )^\top \mathbf{M}_r (\ \mathbf{h} + \mathbf{r} - \mathbf{t}\ )$	$\ \mathbf{h}\ _2 \leq 1, \ \mathbf{t}\ _2 \leq 1, \ \mathbf{r}\ _2 \leq 1$ $\ \mathbf{M}_r\ _F \leq 1, [\mathbf{M}_r]_{ij} = [\mathbf{M}_r]_{ji} \geq 0$
KG2E [45]	$\mathbf{h} \sim \mathcal{N}(\boldsymbol{\mu}_h, \boldsymbol{\Sigma}_h)$ $\mathbf{t} \sim \mathcal{N}(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t)$ $\boldsymbol{\mu}_h, \boldsymbol{\mu}_t \in \mathbb{R}^d$ $\boldsymbol{\Sigma}_h, \boldsymbol{\Sigma}_t \in \mathbb{R}^{d \times d}$	$\mathbf{r} \sim \mathcal{N}(\boldsymbol{\mu}_r, \boldsymbol{\Sigma}_r)$ $\boldsymbol{\mu}_r \in \mathbb{R}^d, \boldsymbol{\Sigma}_r \in \mathbb{R}^{d \times d}$	$-\text{tr}(\boldsymbol{\Sigma}_r^{-1}(\boldsymbol{\Sigma}_h + \boldsymbol{\Sigma}_t)) - \boldsymbol{\mu}^\top \boldsymbol{\Sigma}_r^{-1} \boldsymbol{\mu} - \ln \frac{\det(\boldsymbol{\Sigma}_r)}{\det(\boldsymbol{\Sigma}_h + \boldsymbol{\Sigma}_t)}$ $-\boldsymbol{\mu}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} - \ln(\det(\boldsymbol{\Sigma}))$ $\boldsymbol{\mu} = \boldsymbol{\mu}_h + \boldsymbol{\mu}_r - \boldsymbol{\mu}_t$ $\boldsymbol{\Sigma} = \boldsymbol{\Sigma}_h + \boldsymbol{\Sigma}_r + \boldsymbol{\Sigma}_t$	$\ \boldsymbol{\mu}_h\ _2 \leq 1, \ \boldsymbol{\mu}_t\ _2 \leq 1, \ \boldsymbol{\mu}_r\ _2 \leq 1$ $c_{min} \mathbf{I} \leq \boldsymbol{\Sigma}_h \leq c_{max} \mathbf{I}$ $c_{min} \mathbf{I} \leq \boldsymbol{\Sigma}_t \leq c_{max} \mathbf{I}$ $c_{min} \mathbf{I} \leq \boldsymbol{\Sigma}_r \leq c_{max} \mathbf{I}$
TransG [46]	$\mathbf{h} \sim \mathcal{N}(\boldsymbol{\mu}_h, \sigma_h^2 \mathbf{I})$ $\mathbf{t} \sim \mathcal{N}(\boldsymbol{\mu}_t, \sigma_t^2 \mathbf{I})$ $\boldsymbol{\mu}_h, \boldsymbol{\mu}_t \in \mathbb{R}^d$	$\boldsymbol{\mu}_r^i \sim \mathcal{N}(\boldsymbol{\mu}_t - \boldsymbol{\mu}_h, (\sigma_h^2 + \sigma_t^2) \mathbf{I})$ $\mathbf{r} = \sum_i \pi_r^i \boldsymbol{\mu}_r^i \in \mathbb{R}^d$	$\sum_i \pi_r^i \exp\left(-\frac{\ \boldsymbol{\mu}_h + \boldsymbol{\mu}_r^i - \boldsymbol{\mu}_t\ _2^2}{\sigma_h^2 + \sigma_t^2}\right)$	$\ \boldsymbol{\mu}_h\ _2 \leq 1, \ \boldsymbol{\mu}_t\ _2 \leq 1, \ \boldsymbol{\mu}_r^i\ _2 \leq 1$
UM [56]	$\mathbf{h}, \mathbf{t} \in \mathbb{R}^d$	—	$-\ \mathbf{h} - \mathbf{t}\ _2^2$	$\ \mathbf{h}\ _2 = 1, \ \mathbf{t}\ _2 = 1$
SE [57]	$\mathbf{h}, \mathbf{t} \in \mathbb{R}^d$	$\mathbf{M}_r^1, \mathbf{M}_r^2 \in \mathbb{R}^{d \times d}$	$-\ \mathbf{M}_r^1 \mathbf{h} - \mathbf{M}_r^2 \mathbf{t}\ _1$	$\ \mathbf{h}\ _2 = 1, \ \mathbf{t}\ _2 = 1$

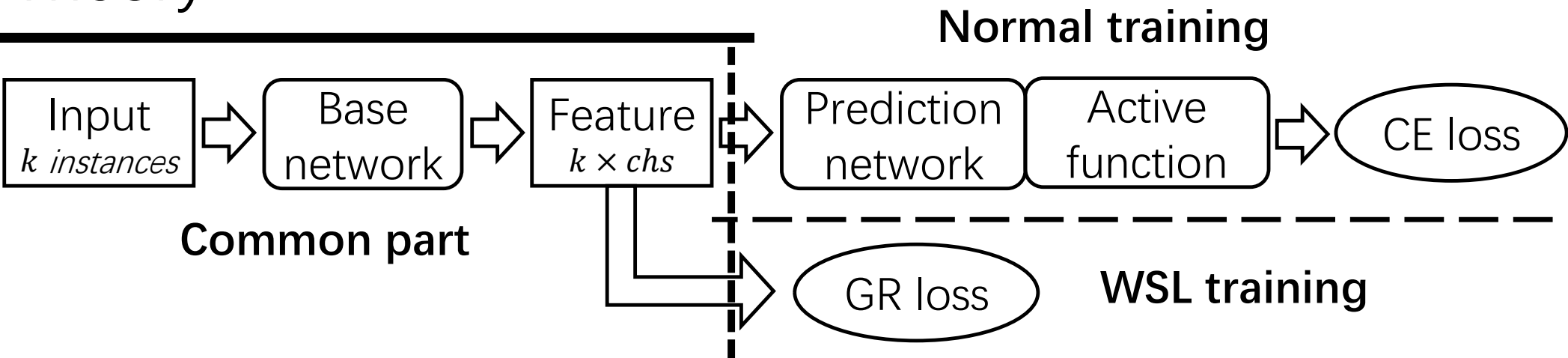
# Weakly Supervised Learning

---

1. Introduction
2. Theory
3. Application and SOTA works
  - Appetizer: image classification task
  - Main-dish: speech recognition task
  - Dessert: what about self-supervised learning
4. Summary



# Theory



Example of **Graph Regularization** loss:  
(loss inside single group)

$$L_{reg} = \sum_{i,j} \tilde{s}_{i,j} \|h(\mathbf{x}_i) - h(\mathbf{x}_j)\|_2^2,$$

Accumulate  
all pairs of  
samples

Weight, calculated  
by A,B, optional.

Instance A,B,  
selected by  
strategy

Distance function, usually L2

What people focus:

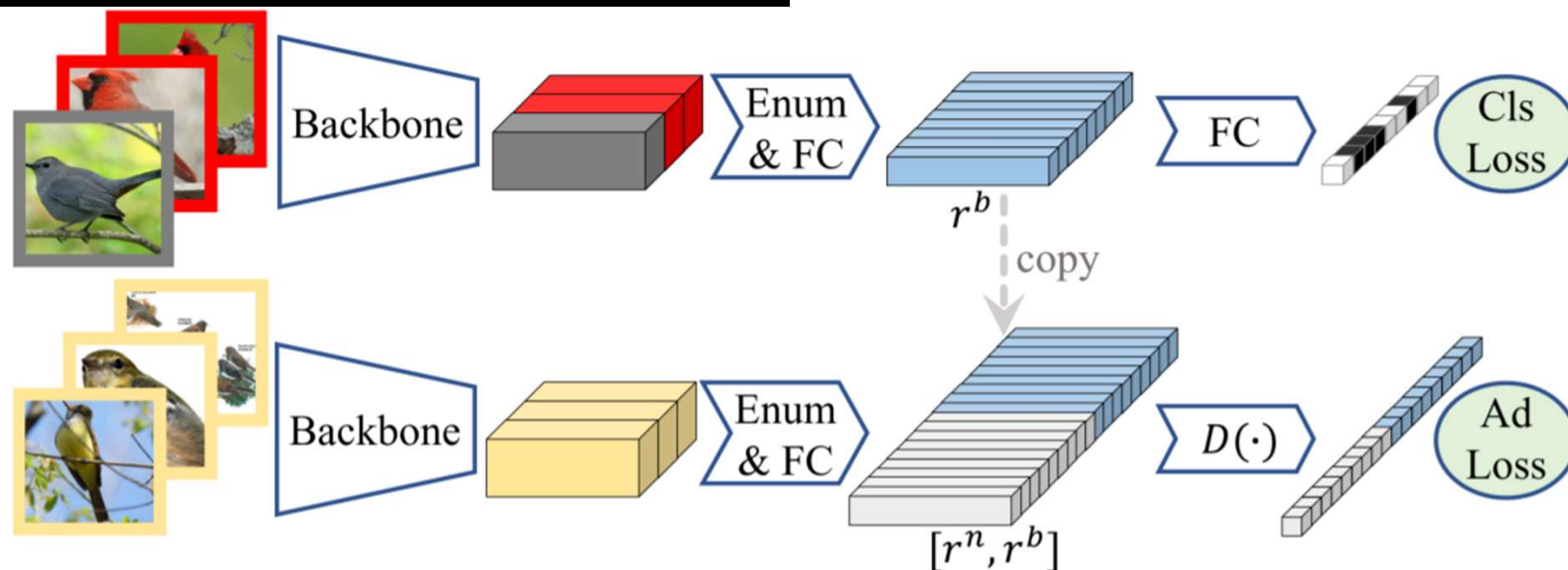
- How to choose instance pairs without label?
- How to add weight on these pairs?
- How to maximize distance between groups?

# Weakly Supervised Learning

---

1. Introduction
2. Theory
3. Application and SOTA works
  - Appetizer: image classification task
  - Main-dish: speech recognition task
  - Dessert: what about self-supervised learning
4. Summary

# Appetizer: Weak-shot image classification



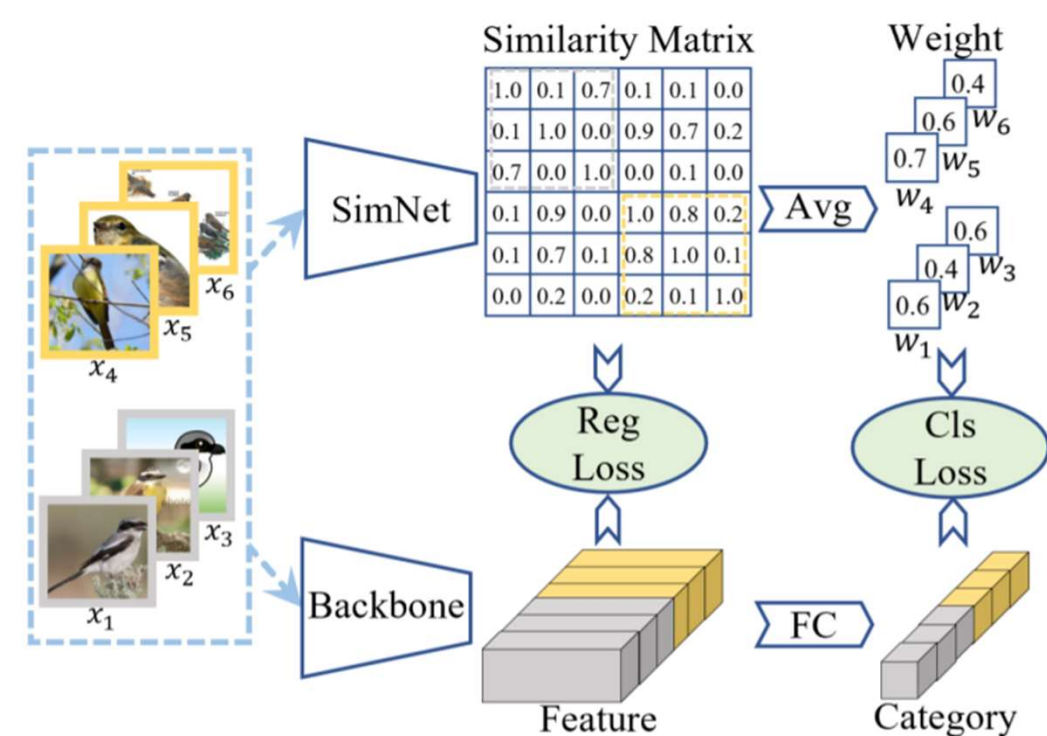
Novelty:

- Use *transferred similarity + denoise strategy* to tackle web training data.
- Apply adversarial loss to similarity net.

- Network (derived from *ref.1.1*):
- Pair-enumerate:  $(k, d) \rightarrow (k, k, 2d) \rightarrow (k^2, 2d)$
- Fc:  $(k^2, 2d) \rightarrow (k^2)$ ; Similar prediction: Y/N
- .....
- Feature constraint: cluster unlabeled data based
- on feature from labeled data.

*Weak-shot Fine-grained Classification via Similarity Transfer, NIPS 2021*

# Appetizer: Weak-shot image classification



Training steps:

1. Train network on labeled data.
  - Prediction target: binary classification of similar or not.
  - CE loss.
  - Strong constraint: batch size  $k$  s.t.  $k^2$  items pre-step.
2. Cluster unlabeled data via pre-trained network.
  - Select  $k_1$  samples from labeled data.
  - Select  $k_2$  samples from unlabeled data.
  - Calculate feature  $f_1, f_2$ .
  - Weighted L2 loss by similarity.

$$L_{reg} = \sum_{i,j} \tilde{s}_{i,j} \|h(\mathbf{x}_i) - h(\mathbf{x}_j)\|_2^2, \quad w_{c,i} = \frac{1}{N_c^n} \sum_{j=1}^{N_c^n} \frac{s_{c,i,j} + s_{c,j,i}}{2}.$$

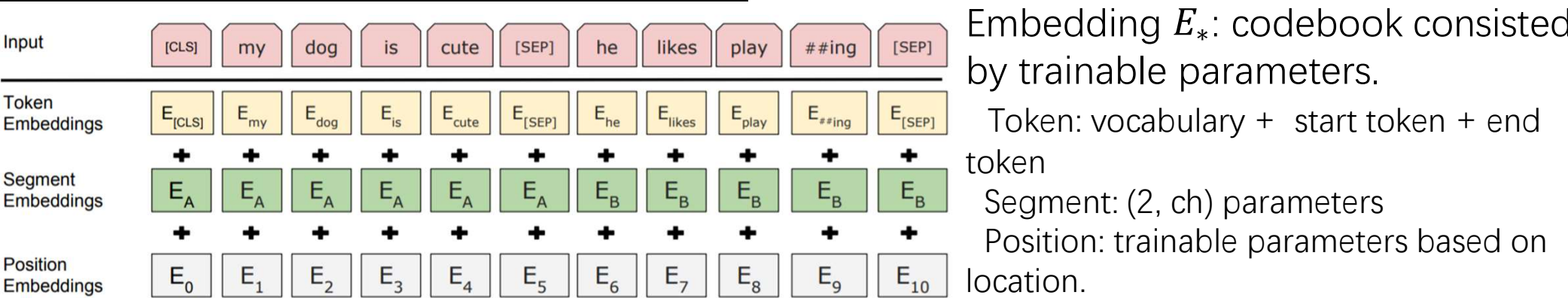
*Weak-shot Fine-grained Classification via Similarity Transfer, NIPS 2021*

# Weakly Supervised Learning

---

1. Introduction
2. Theory
3. Application and SOTA works
  - Appetizer: image classification task
  - Main-dish: speech recognition task
  - Dessert: what about self-supervised learning
4. Summary

# Main-dish: W2v-BERT



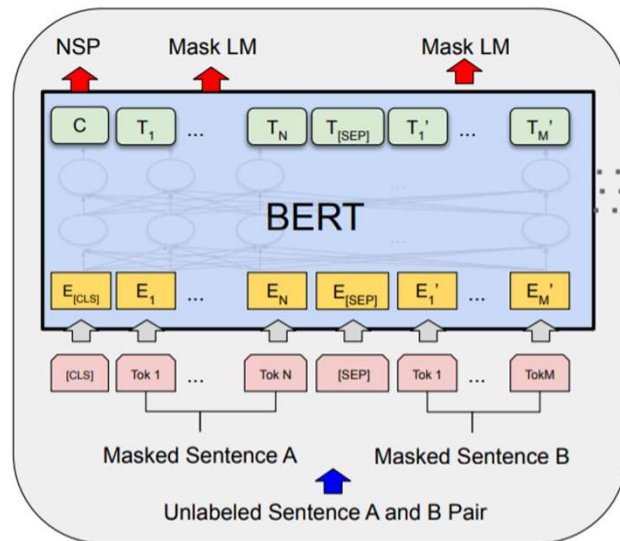
Input: 2 sentences A,B from articles.

Task: randomly mask multi-words (15%) in A and B, predict the words by context (CE loss).

Task: randomly select continuous and discontinuous (50%/50%) sentences A and B, determine whether B is the next sentence of A (CE loss).

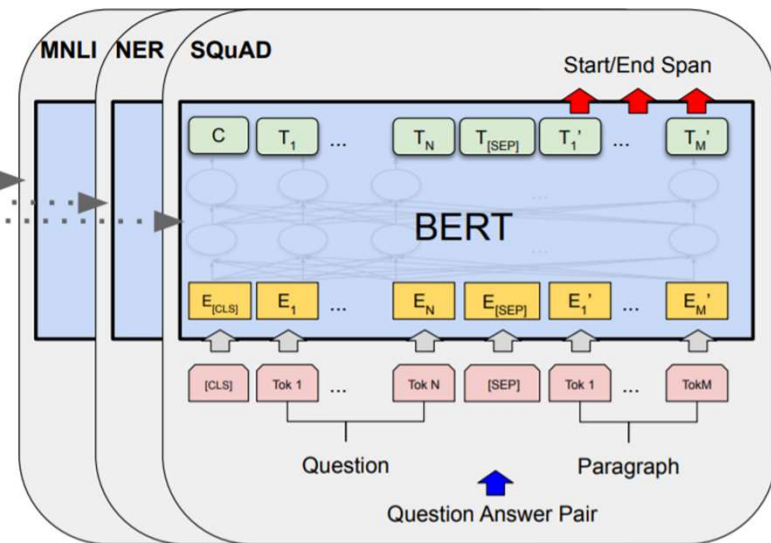
*BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding, NACCL 2019, Google AI*

# Main-dish: W2v-BERT



Pre-training

Bidirectional encoder network: transformer with unfixed sequence.

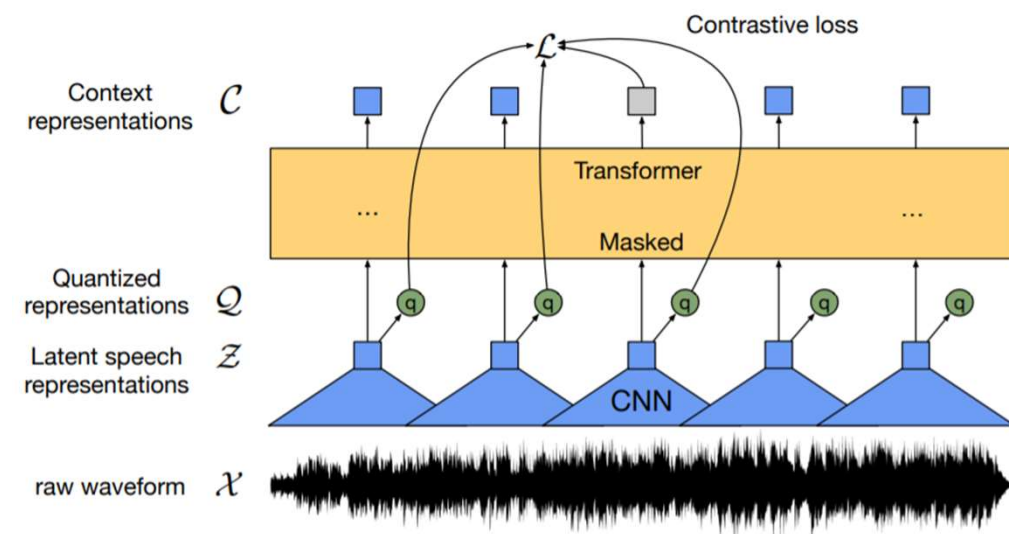


Fine-Tuning

Normal supervised training for task with low amount (e.g. 1/1000 of unlabeled data), labeled, specific data (so called downstream task).

*BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding, NACCL 2019, Google AI*

# Main-dish: W2v-BERT



Auxiliary task: identify the true quantized representation from masked region (*ref2.1.1*).

$G$  codebooks  $\in \mathbb{R}^{V \times d}$ , product quantization:

$z \in \mathbb{R}^{ch} \xrightarrow{f(z)} I \in \mathbb{R}^{G \times V} \xrightarrow{\text{Gum.Smax}(I)} p_{g,v}$

select  $G$  rows  $e_1, \dots, e_G$  by  $p_{g,v}$  via *argmax*

concatenate into  $e \in \mathbb{R}^{G \times d} \xrightarrow{f(z)} q \in \mathbb{R}^{ch}$

Gumbel softmax:  $p_{g,v} = \frac{\exp(l_{g,v} + n_v)/\tau}{\sum_{k=1}^V \exp(l_{g,k} + n_k)/\tau}$ ,

$\tau$ : non-negative temperature

$n = -\log(-\log(u))$

$u$ : uniform sample from (0,1)

Input: speech voice (1d sequence, float)

Task 1: identify the true quantized latent speech representation

Task 2: keep vectors in codebooks used as equal as possible.

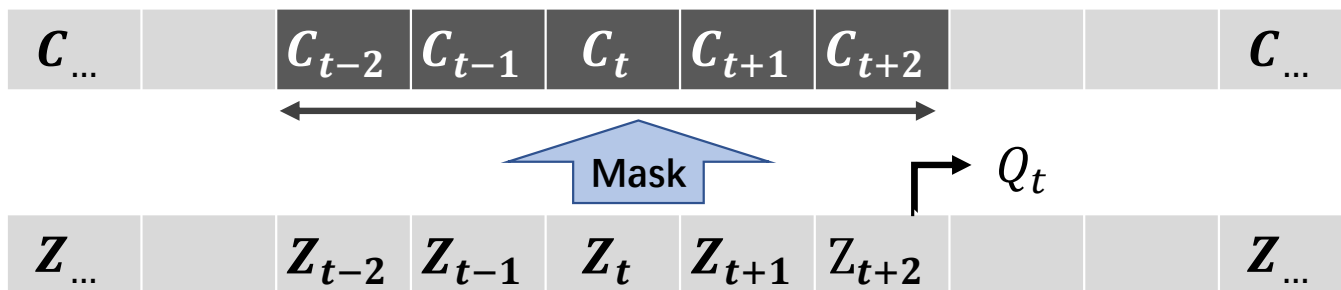
*Wav2vec 2.0: A Framework for Self-Supervised Learning of Speech Representations, NIPS 2020, Facebook AI*



# Main-dish: W2v-BERT

Random mask:

replace inputs of transformer by shared, trained feature vector



Contrastive Loss (*ref2.1.2*):

$$\mathcal{L}_m = -\log \frac{\exp(\text{sim}(\mathbf{c}_t, \mathbf{q}_t)/\kappa)}{\sum_{\tilde{\mathbf{q}} \sim \mathbf{Q}_t} \exp(\text{sim}(\mathbf{c}_t, \tilde{\mathbf{q}})/\kappa)}$$

$\mathbf{c}$ : context representation

$\mathbf{q}$ : quantized representation

$\text{sim}$ : cosine similarity

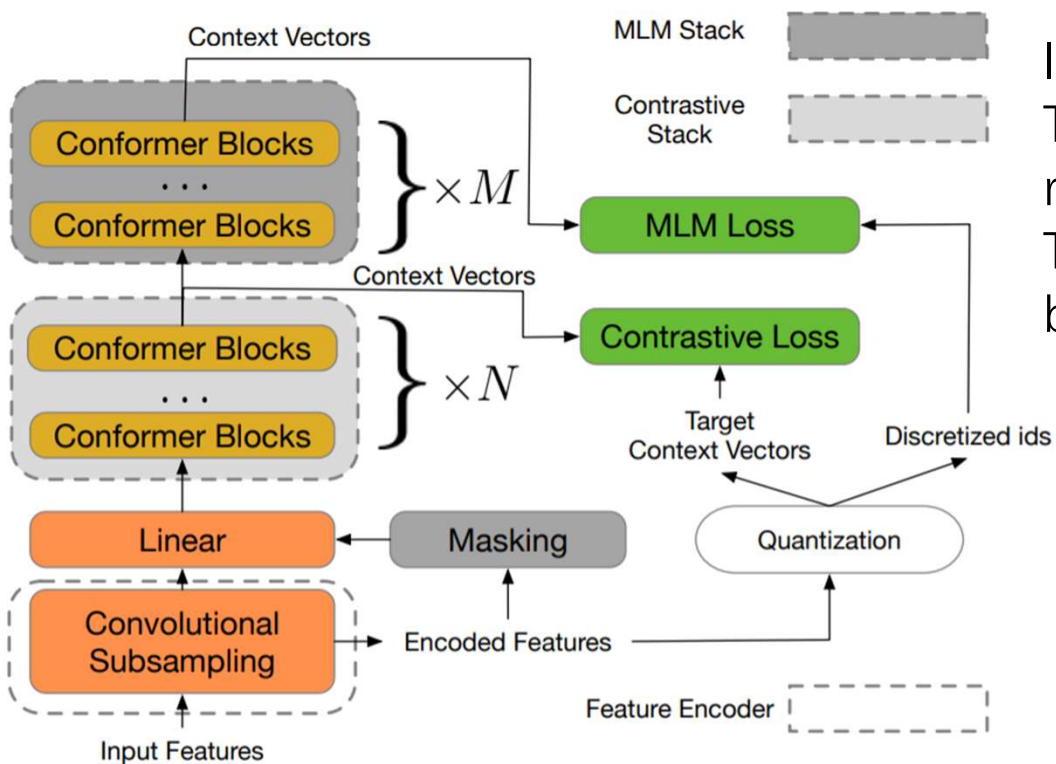
$$\text{sim}(\mathbf{a}, \mathbf{b}) = \mathbf{a}^T \mathbf{b} / \|\mathbf{a}\| \|\mathbf{b}\|$$

$$\mathcal{L}_d = \frac{1}{GV} \sum_{g=1}^G -H(\bar{p}_g) = \frac{1}{GV} \sum_{g=1}^G \sum_{v=1}^V \bar{p}_{g,v} \log \bar{p}_{g,v}$$

Diversity Loss: encourage the equal use of the  $V$  entries in each of the  $G$  codebooks by maximizing the entropy of the averaged softmax distribution  $I$

*Wav2vec 2.0: A Framework for Self-Supervised Learning of Speech Representations, NIPS 2020, Facebook AI*

# Main-dish: W2v-BERT



Combine W2v and BERT to conduct WSL on unlabeled voice data.

Input: unlabeled voice data

Task: identify the true quantized latent speech representation (contrastive loss)

Task: randomly mask words and predict them by context (**M**asked **L**anguage **M**odeling).

Contrastive stack: transformer-based encoder

MLM stack: transformer-based decoder

*W2v-BERT: Combining Contrastive Learning and Masked Language Modeling for Self-Supervised Speech Pre-Training, 2021, MIT & Google Brain*

# Main-dish: W2v-BERT

Method	Unlabeled Data (hrs)	AM Size (B)	LM Size (B)	No LM				With LM			
				dev	dev-other	test	test-other	dev	dev-other	test	test-other
<b>Trained from Scratch</b>											
Conformer L [21]*	N/A	0.1	0.1	1.9	4.4	2.1	4.3	—	—	1.9	3.9
<b>Self-training Only</b>											
Conformer L with NST [21]	60k	0.1	0.1	1.6	3.3	1.7	3.5	1.6	3.1	1.7	3.3
<b>Pre-training Only</b>											
wav2vec 2.0 [22]	60k	0.3	> 0.4 <sup>†</sup>	2.1	4.5	2.2	4.5	1.6	3.0	1.8	3.3
HuBERT Large [25]	60k	0.3	—	—	—	—	—	1.5	3.0	1.9	3.3
HuBERT X-Large [25]	60k	1.0	—	—	—	—	—	1.5	<b>2.5</b>	1.8	2.9
w2v-Conformer XL [21]	60k	0.6	0.1	1.7	3.5	1.7	3.5	1.6	3.2	<b>1.5</b>	3.2
w2v-Conformer XXL [21]	60k	1.0	0.1	1.6	3.2	1.6	3.3	1.5	3.0	<b>1.5</b>	3.1
w2v-BERT XL (Ours)	60k	0.6	0.1	<b>1.5</b>	2.9	<b>1.5</b>	2.9	<b>1.4</b>	2.8	<b>1.5</b>	2.8
w2v-BERT XXL (Ours)	60k	1.0	0.1	<b>1.5</b>	<b>2.7</b>	<b>1.5</b>	<b>2.8</b>	<b>1.4</b>	2.6	<b>1.5</b>	<b>2.7</b>
<b>Pre-training + Self-training</b>											
wav2vec 2.0 [22]	60k	0.3	> 0.4	<b>1.3</b>	3.1	1.7	3.5	<b>1.1</b>	2.7	1.5	3.1
w2v-Conformer XXL [21]	60k	1.0	0.1	<b>1.3</b>	2.7	1.5	2.8	1.3	2.6	<b>1.4</b>	2.7
w2v-Conformer XXL+ [21]	60k	1.1	0.1	<b>1.3</b>	2.7	1.5	2.7	1.3	2.6	<b>1.4</b>	2.6
w2v-BERT XL (Ours)	60k	0.6	0.1	<b>1.3</b>	2.6	<b>1.4</b>	2.7	1.3	2.6	<b>1.4</b>	2.6
w2v-BERT XXL (Ours)	60k	1.0	0.1	1.4	<b>2.4</b>	<b>1.4</b>	<b>2.5</b>	1.3	<b>2.4</b>	<b>1.4</b>	<b>2.5</b>

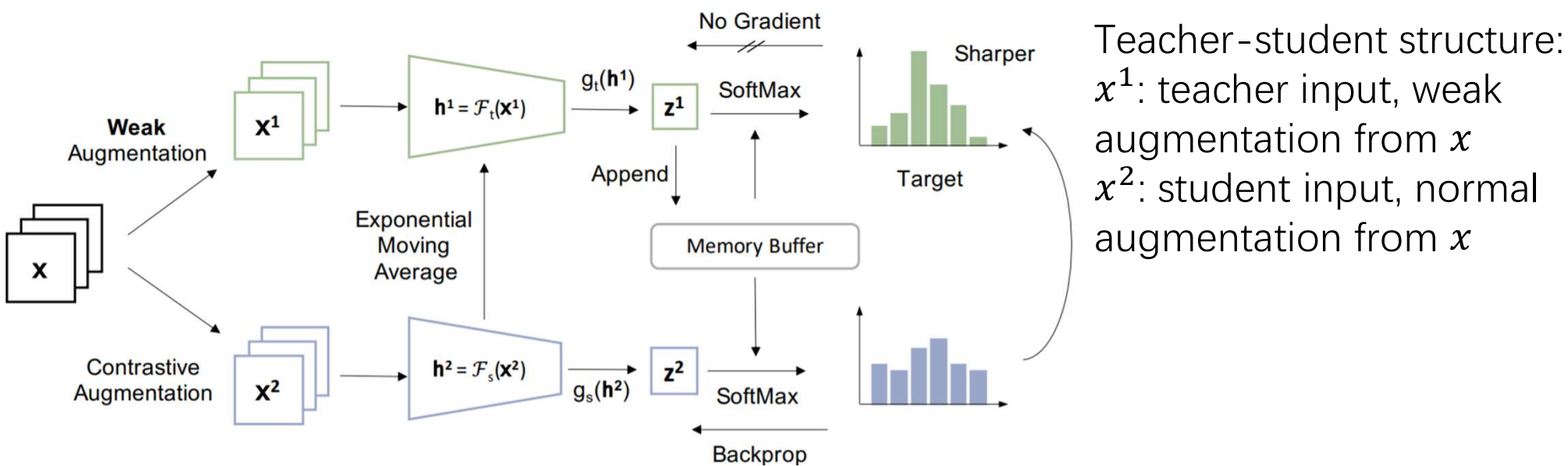
*W2v-BERT: Combining Contrastive Learning and Masked Language Modeling for Self-Supervised Speech Pre-Training, 2021, MIT & Google Brain*

# Weakly Supervised Learning

---

1. Introduction
2. Theory
3. Application and SOTA works
  - Appetizer: image classification task
  - Main-dish: speech recognition task
  - Dessert: what about self-supervised learning
4. Summary

# Dessert: SSL, strategy purely for pre-training



Data augmentation based **Self-Supervised Learning**: distribution of predicted classes of instances between two augmentations should be similar.  
Calculate the similarity between *single*  $z_i^1$  and *all*  $z_j^2$ .

*ReSSL: Relational Self-Supervised Learning with Weak Augmentation, NIPS 2021*

# Dessert: SSL, strategy purely for pre-training

Preliminaries:

## Noise Contrastive Estimation

$$\mathcal{L}_{NCE} = -\log \frac{\exp(\text{sim}(\mathbf{z}^1, \mathbf{z}^2)/\tau)}{\exp(\text{sim}(\mathbf{z}_i^1, \mathbf{z}_i^2)/\tau) + \sum_{k=1}^N \exp(\text{sim}(\mathbf{z}_i^1, \mathbf{z}_k)/\tau)}$$

Reference (1) points to  $\mathbf{z}^1$  in the numerator.  
 Non-zero  $\tau$  points to  $\tau$  in the numerator.  
 Target term (1) points to  $\text{sim}(\mathbf{z}_i^1, \mathbf{z}_i^2)/\tau$  in the denominator.  
 Local2global term ( $1 \times N$ ) points to  $\sum_{k=1}^N \exp(\text{sim}(\mathbf{z}_i^1, \mathbf{z}_k)/\tau)$  in the denominator.

*sim*: cosine similarity

$$\text{sim}(\mathbf{u}, \mathbf{v}) = \mathbf{u}^T \mathbf{v} / \|\mathbf{u}\| \|\mathbf{v}\|$$

Weak aug. and normal aug. vs. original image:

$$\mathbf{p}_i^1 = \frac{\exp(\text{sim}(\mathbf{z}^1, \mathbf{z}_i)/\tau_t)}{\sum_{k=1}^K \exp(\text{sim}(\mathbf{z}^1, \mathbf{z}_k)/\tau_t)}, \quad \mathbf{p}_i^2 = \frac{\exp(\text{sim}(\mathbf{z}^2, \mathbf{z}_i)/\tau_s)}{\sum_{k=1}^K \exp(\text{sim}(\mathbf{z}^2, \mathbf{z}_k)/\tau_s)}$$

1-1 similarity (top arrow)  
 1-all similarity (bottom arrow)

$\mathbf{z}^{1,2}$ : 1<sup>st</sup> image from weak and normal aug.

$\mathbf{z}_i$ :  $i$ -th original image in batch.

Minimize Kullback–Leibler divergence

$$\mathcal{L}_{relation} = D_{KL}(\mathbf{p}^1 || \mathbf{p}^2) = H(\mathbf{p}^1, \mathbf{p}^2) - H(\mathbf{p}^1)$$

Because  $H(\mathbf{p}^1)$  is target distribution so only regress  $H(\mathbf{p}^1, \mathbf{p}^2)$  hence  $L = H_{CE}(\mathbf{p}^1, \mathbf{p}^2)$

*ReSSL: Relational Self-Supervised Learning with Weak Augmentation, NIPS 2021*

# Weakly Supervised Learning

---

1. Introduction
2. Theory
3. Application and SOTA works
  - Appetizer: image classification task
  - Main-dish: speech recognition task
  - Dessert: what about self-supervised learning
4. Summary

# WSL: a way to utilize public, unlabeled data

---

## Methodology:

- Clustering on feature space: reduce feature distance, increase similarity, etc.
- Data augmentation: add uncertainty on certain data.

## From view of business:

- Convert pre-trained network as asset: one network to multiple downstream tasks.
- Reduce labeling cost.
- Reduce data storage cost: unlabeled data can be removed after producing pre-trained network.

## From view of R&D:

- Additional supervision is conducive to regress network more quickly, precisely.
- Graph based supervision positively contribute to class robustness.
- Flexible, dynamic network structure to deliver: FC layer is no longer the only choice.
- Easily used as novelty in research.