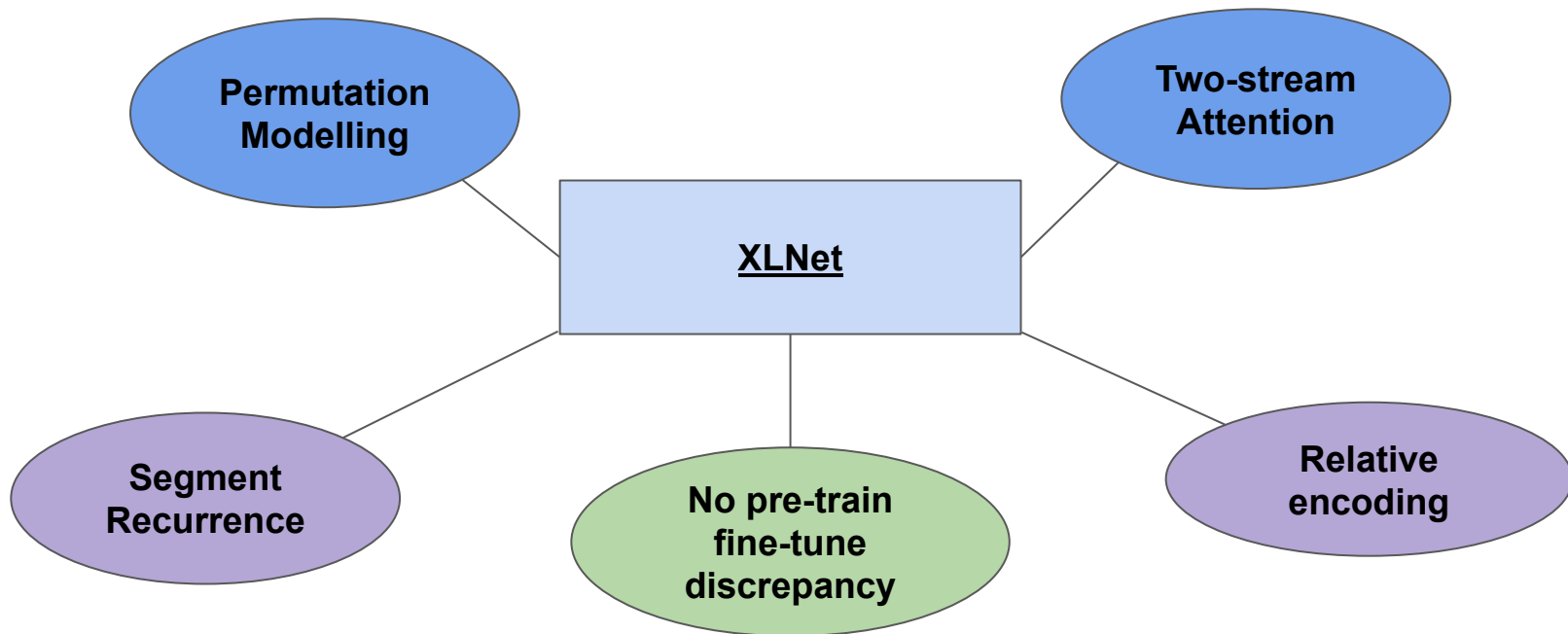# XLNet

Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Ruslan Salakhutdinov, Quoc V. Le
*Team members: Shruti Jha, Shubham Agrawal, Siddhant Agarwal, Simran, Siruvuri Karan Raju, Siya Garg*
**Group No: 20**

## Key Contributions

## Autoregressive Modelling:

$$\max_{\theta} \quad \log p_\theta(\mathbf{x}) = \sum_{t=1}^{T} \log p_\theta(x_t \mid \mathbf{x}_{<t})$$

## Autoencoding Modelling:

$$\max_{\theta} \quad \log p_\theta(\bar{\mathbf{x}} \mid \hat{\mathbf{x}}) \approx \sum_{t=1}^{T} m_t \log p_\theta(x_t \mid \hat{\mathbf{x}})$$

$$\sum_{t=1}^{T} \log p_\theta(x_t \mid \mathbf{x}_{<t}) = \sum_{t=1}^{T} \log \frac{\exp\left(h_\theta(\mathbf{x}_{1:t-1})^\top e(x_t)\right)}{\sum_{x'} \exp\left(h_\theta(\mathbf{x}_{1:t-1})^\top e(x')\right)}$$

$$\sum_{t=1}^{T} m_t \log p_\theta(x_t \mid \hat{\mathbf{x}}) = \sum_{t=1}^{T} m_t \log \frac{\exp\left(H_\theta(\hat{\mathbf{x}})_t^\top e(x_t)\right)}{\sum_{x'} \exp\left(H_\theta(\hat{\mathbf{x}})_t^\top e(x')\right)}$$

pretrain-finetune discrepancy



Factorization order: 3 → 2 → 4 → 1

Factorization order: 2 → 4 → 3 → 1

Factorization order: 1 → 4 → 2 → 3

Factorization order: 4 → 3 → 1 → 2

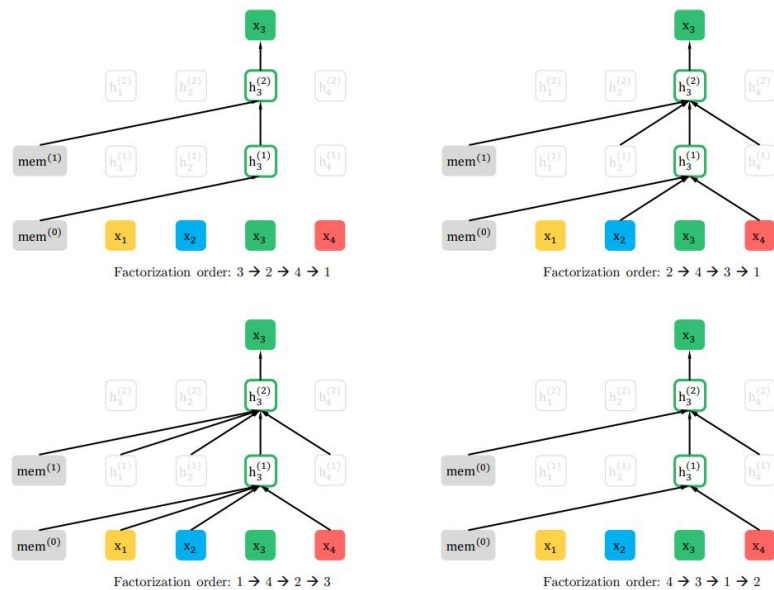**Training Regime:** Autoregressive

$$\mathcal{J}_{\text{BERT}} = \log p(\text{New} \mid \text{is a city}) + \log p(\text{York} \mid \text{is a city}),$$

$$\mathcal{J}_{\text{XLNet}} = \log p(\text{New} \mid \text{is a city}) + \log p(\text{York} \mid \text{New}, \text{is a city}).$$

## Two-Stream Self-Attention for Target-Aware Representations

$$g_{z_t}^{(m)} \leftarrow \text{Attention}(\mathbf{Q} = g_{z_t}^{(m-1)}, \mathbf{KV} = \mathbf{h}_{\mathbf{z}_{<t}}^{(m-1)}; \theta), \quad \text{(query stream: use } z_t \text{ but cannot see } x_{z_t})$$

$$h_{z_t}^{(m)} \leftarrow \text{Attention}(\mathbf{Q} = h_{z_t}^{(m-1)}, \mathbf{KV} = \mathbf{h}_{\mathbf{z}_{\leq t}}^{(m-1)}; \theta), \quad \text{(content stream: use both } z_t \text{ and } x_{z_t}).$$

**OPTIMIZATION:** Partial Prediction

Introduced in XLNet

Standard Transformer

# Performance/Model Comparison

| Model | MNLI | QNLI | QQP | RTE | SST-2 | MRPC | CoLA | STS-B | WNLI |
|---|---|---|---|---|---|---|---|---|---|
| *Single-task single models on dev* | | | | | | | | | |
| BERT [2] | 86.6/- | 92.3 | 91.3 | 70.4 | 93.2 | 88.0 | 60.6 | 90.0 | - |
| RoBERTa [21] | 90.2/90.2 | 94.7 | 92.2 | **86.6** | 96.4 | **90.9** | 68.0 | 92.4 | - |
| XLNet | **90.8/90.8** | **94.9** | **92.3** | 85.9 | **97.0** | 90.8 | **69.0** | **92.5** | - |
| *Multi-task ensembles on test (from leaderboard as of Oct 28, 2019)* | | | | | | | | | |
| MT-DNN* [20] | 87.9/87.4 | 96.0 | 89.9 | 86.3 | 96.5 | 92.7 | 68.4 | 91.1 | 89.0 |
| RoBERTa* [21] | 90.8/90.2 | 98.9 | 90.2 | 88.2 | 96.7 | 92.3 | 67.8 | 92.2 | 89.0 |
| XLNet* | **90.9/90.9**$^\dagger$ | **99.0**$^\dagger$ | **90.4**$^\dagger$ | **88.5** | **97.1**$^\dagger$ | **92.9** | **70.2** | **93.0** | **92.5** |

| # | Model | RACE | SQuAD2.0 | | MNLI | SST-2 |
|---|---|---|---|---|---|---|
| | | | F1 | EM | m/mm | |
| 1 | BERT-Base | 64.3 | 76.30 | 73.66 | 84.34/84.65 | 92.78 |
| 2 | DAE + Transformer-XL | 65.03 | 79.56 | 76.80 | 84.88/84.45 | 92.60 |
| 3 | XLNet-Base ($K = 7$) | 66.05 | **81.33** | **78.46** | **85.84/85.43** | 92.66 |
| 4 | XLNet-Base ($K = 6$) | 66.66 | 80.98 | 78.18 | 85.63/85.12 | **93.35** |
| 5 | - memory | 65.55 | 80.15 | 77.27 | 85.32/85.05 | 92.78 |
| 6 | - span-based pred | 65.95 | 80.61 | 77.91 | 85.49/85.02 | 93.12 |
| 7 | - bidirectional data | 66.34 | 80.65 | 77.87 | 85.31/84.99 | 92.66 |
| 8 | + next-sent pred | **66.76** | 79.83 | 76.94 | 85.32/85.09 | 92.89 |