

Tutorial on Transformers with Indices

Ref: Section I-0, “Crash Course on Transformers” (5 pages)
of [[arXiv:2304.02034](https://arxiv.org/abs/2304.02034)] (Emily Dinan, Sho Yaida, Susan Zhang)

[cf. Andrej Karpathy’s YouTube tutorial]

Outline

0. Preliminary: Vocabulary and Task
1. Word Embedding and Positional Embedding
2. Layer Normalizations and Skip Connections
3. MLP Blocks
4. Multi-Head Self-Attention Blocks
5. Head

0. Preliminary: Vocabulary and Task

Example 1

Vocabulary=a set of characters

$$\{a, b, \dots, !, ?, [\text{space}]\}$$

0. Preliminary: Vocabulary and Task

Example 1

Vocabulary=a set of characters

$$\{a, b, \dots, !, ?, [\text{space}]\}$$

$0 \quad 1 \quad \quad \quad n_{\text{vocab}} - 1$

Task=next-character prediction

Input: B|r|e|a|k|i|n|g|[space]|B|a

Output: B|r|e|a|k|i|n|g|[space]|B|a|d

0. Preliminary: Vocabulary and Task

Example 2

Vocabulary=a set of “tokens” (created by tokenization such as byte-pair encoding algorithm)

{the, of, ..., super, ..., symmetry, ..., [space]}

Input Notation

$$\mathcal{X}_{\alpha;t;i}$$

sample 1: Super|symmetry|[space]|is|[space]|awesome|.

sample 2: Transformers|[space]|are|[space]|simple|.

sample 3: I|[space]|want|[space]|to|[space]|watch|[space]|Breaking|[space]|Bad|.

...

Input Notation

$$x_{\alpha;t;i}$$

$\alpha = 1, \dots, N_{\text{sample}}$

sample 1: Super|symmetry|[space]|is|[space]|awesome|.

sample 2: Transformers|[space]|are|[space]|simple|.

sample 3: I|[space]|want|[space]|to|[space]|watch|[space]|Breaking|[space]|Bad|.

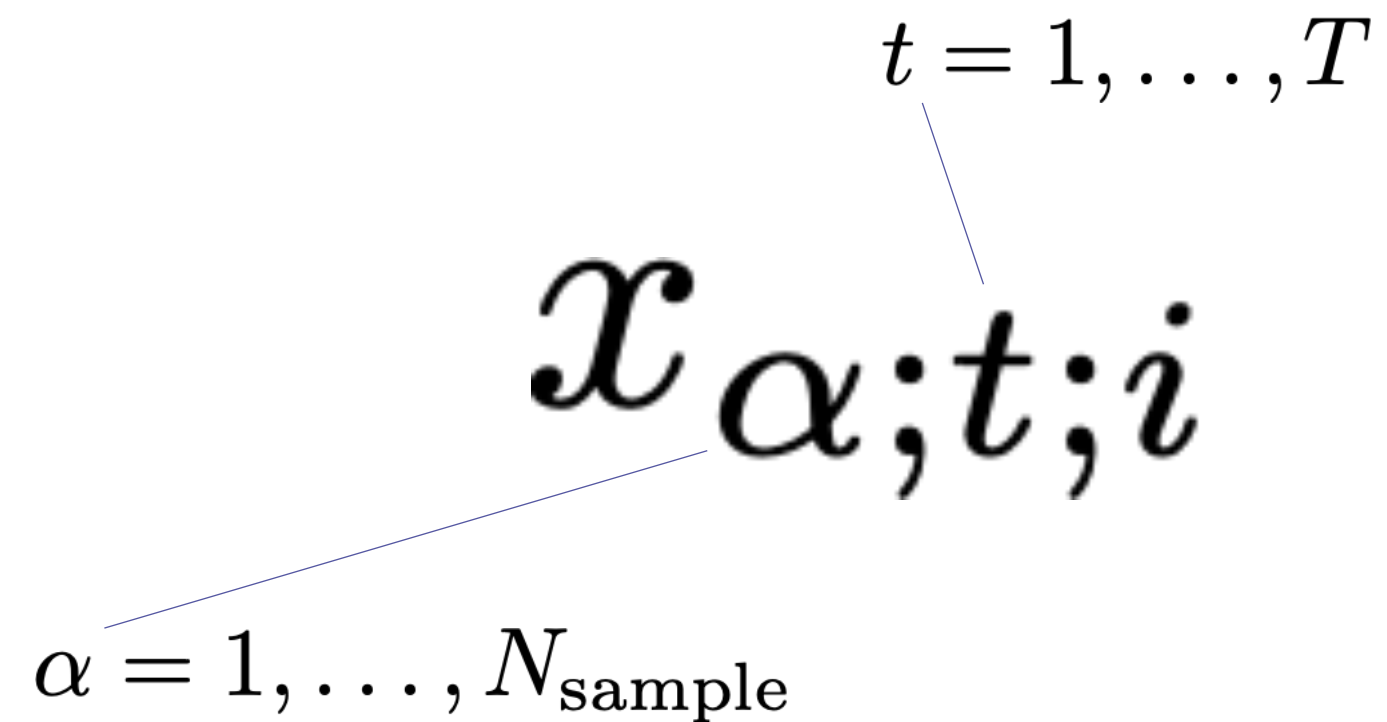
...

Input Notation

$$x_{\alpha;t;i}$$

$t = 1, \dots, T$

$\alpha = 1, \dots, N_{\text{sample}}$



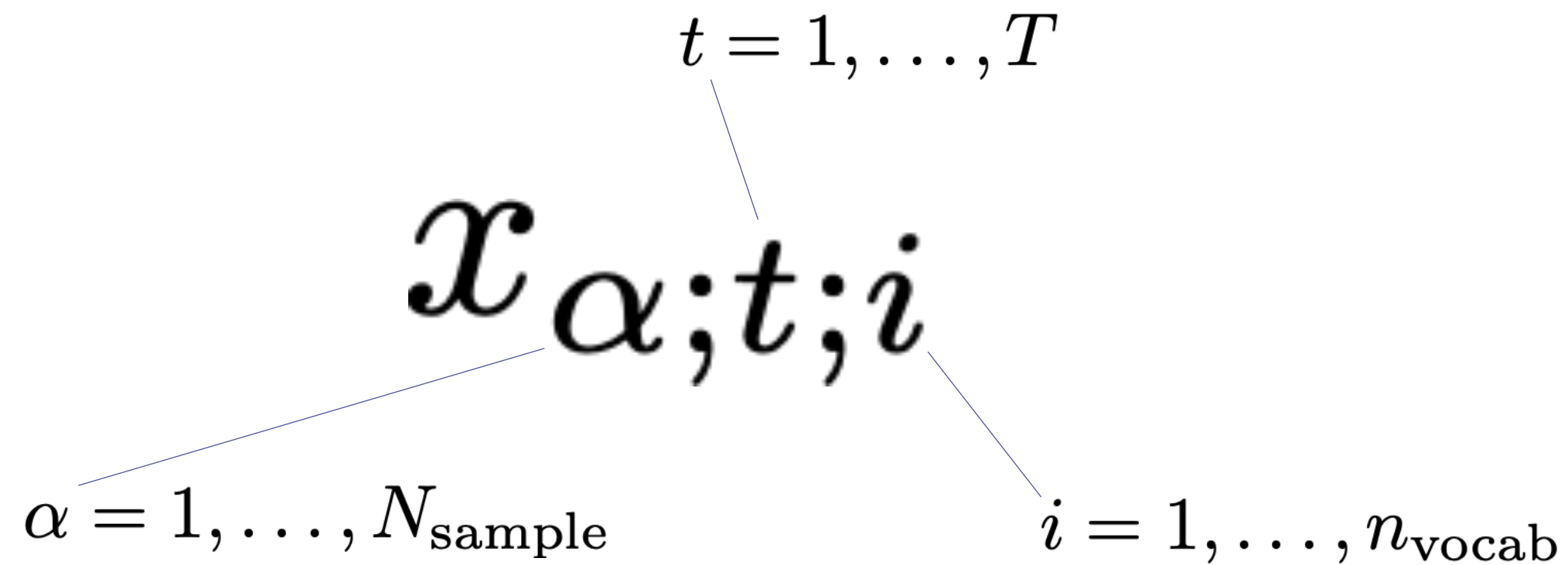
sample 1: Super|symmetry|[space]|is|[space]|awesome|.

sample 2: Transformers|[space]|are|[space]|simple|.

sample 3: |[space]|want|[space]|to|[space]|watch|[space]|Breaking|[space]|Bad|.

...

Input Notation



one-hot in i -direction: $(0, \dots, 0, 1, 0, \dots, 0)$

sample 1: Super|symmetry|[space]|is|[space]|awesome|.

sample 2: Transformers|[space]|are|[space]|simple|.

sample 3: |[space]|want|[space]|to|[space]|watch|[space]|Breaking|[space]|Bad|.

...

1. Word Embedding and Positional Embedding

$$\tilde{z}_{\alpha;t;i} = \sum_{j=1}^{n_{\text{vocab}}} W_{ij}^{\text{WE}} x_{\alpha;t;j} = W_{ij^*(\alpha;t)}^{\text{WE}} \quad i = 1, \dots, n$$

1. Word Embedding and Positional Embedding

$$\tilde{z}_{\alpha;t;i} = \sum_{j=1}^{n_{\text{vocab}}} W_{ij}^{\text{WE}} x_{\alpha;t;j} = W_{ij^*(\alpha;t)}^{\text{WE}} \quad i = 1, \dots, n$$

$$z_{\alpha;t;i}^{(1)} = b_{t;i}^{\text{PE}} + \tilde{z}_{\alpha;t;i}$$

2. Layer Normalizations and Skip Connections

$$s_{\alpha;t;i}^{(\ell)} = \gamma_i^{(\ell)} \left[\frac{z_{\alpha;t;i}^{(\ell)} - \left(\frac{1}{n} \sum_{j=1}^n z_{\alpha;t;j}^{(\ell)} \right)}{\sqrt{\frac{1}{n} \sum_{j=1}^n \left(z_{\alpha;t;j}^{(\ell)} \right)^2 - \left(\frac{1}{n} \sum_{j=1}^n z_{\alpha;t;j}^{(\ell)} \right)^2 + \epsilon}} \right] + \beta_i^{(\ell)}$$

2. Layer Normalizations and Skip Connections

$$s_{\alpha;t;i}^{(\ell)} = \gamma_i^{(\ell)} \left[\frac{z_{\alpha;t;i}^{(\ell)} - \left(\frac{1}{n} \sum_{j=1}^n z_{\alpha;t;j}^{(\ell)} \right)}{\sqrt{\frac{1}{n} \sum_{j=1}^n \left(z_{\alpha;t;j}^{(\ell)} \right)^2 - \left(\frac{1}{n} \sum_{j=1}^n z_{\alpha;t;j}^{(\ell)} \right)^2 + \epsilon}} \right] + \beta_i^{(\ell)}$$

$$z_{\alpha;t;i}^{(\ell+1)} = \mathbf{R}_{t;i}^{(\ell+1)} \left(s_{\alpha}^{(\ell)}; \theta^{(\ell+1)} \right) + z_{\alpha;t;i}^{(\ell)}$$

3. MLP Blocks $\mathbf{R}_{t;i}^{\text{MLP}}(s_\alpha; \theta)$

$$w_{\alpha;t;i} = \sum_{j=1}^n W_{ij} s_{\alpha;t;j} \quad \text{for } i = 1, \dots, Mn,$$

$$r_{\alpha;t;i} = \sum_{j=1}^{Mn} X_{ij} \sigma(w_{\alpha;t;j}) \quad \text{for } i = 1, \dots, n,$$

4. Multi-Head Self-Attention Blocks $\mathbf{R}_{t;i}^{\text{MHSA}}(s_\alpha; \theta)$

$$q_{\alpha;t;c}^h \equiv \sum_{i=1}^n Q_{ci}^h s_{\alpha;t;i},$$

$$t = 1, \dots, T$$

$$h = 1, \dots, H$$

$$k_{\alpha;t;c}^h \equiv \sum_{i=1}^n K_{ci}^h s_{\alpha;t;i},$$

$$c = 1, \dots, C \quad \text{where} \quad C = \frac{n}{H}$$

$$v_{\alpha;t;c}^h \equiv \sum_{i=1}^n V_{ci}^h s_{\alpha;t;i},$$

4. Multi-Head Self-Attention Blocks $\mathbf{R}_{t;i}^{\text{MHSA}}(s_\alpha; \theta)$

$$q_{\alpha;t;c}^h \equiv \sum_{i=1}^n Q_{ci}^h s_{\alpha;t;i},$$

$$t = 1, \dots, T$$

$$h = 1, \dots, H$$

$$k_{\alpha;t;c}^h \equiv \sum_{i=1}^n K_{ci}^h s_{\alpha;t;i},$$

$$c = 1, \dots, C \quad \text{where} \quad C = \frac{n}{H}$$

$$v_{\alpha;t;c}^h \equiv \sum_{i=1}^n V_{ci}^h s_{\alpha;t;i},$$

$$\tilde{\Omega}_{\alpha;tt'}^h \equiv \frac{1}{\sqrt{C}} \sum_{c=1}^C q_{\alpha;t;c}^h k_{\alpha;t';c}^h$$

4. Multi-Head Self-Attention Blocks

$$\mathbf{R}_{t;i}^{\text{MHSA}}(s_\alpha; \theta)$$

$$q_{\alpha;t;c}^h \equiv \sum_{i=1}^n Q_{ci}^h s_{\alpha;t;i},$$

$$k_{\alpha;t;c}^h \equiv \sum_{i=1}^n K_{ci}^h s_{\alpha;t;i},$$

$$v_{\alpha;t;c}^h \equiv \sum_{i=1}^n V_{ci}^h s_{\alpha;t;i},$$

$$\tilde{\Omega}_{\alpha;tt'}^h \equiv \frac{1}{\sqrt{C}} \sum_{c=1}^C q_{\alpha;t;c}^h k_{\alpha;t';c}^h$$

Encoder:

$$\Omega_{\alpha;tt'}^h \equiv \frac{\exp(\tilde{\Omega}_{\alpha;tt'}^h)}{\sum_{t''=1}^T \exp(\tilde{\Omega}_{\alpha;tt''}^h)}$$

4. Multi-Head Self-Attention Blocks

$$\mathbf{R}_{t;i}^{\text{MHSA}}(s_\alpha; \theta)$$

$$q_{\alpha;t;c}^h \equiv \sum_{i=1}^n Q_{ci}^h s_{\alpha;t;i},$$

$$k_{\alpha;t;c}^h \equiv \sum_{i=1}^n K_{ci}^h s_{\alpha;t;i},$$

$$v_{\alpha;t;c}^h \equiv \sum_{i=1}^n V_{ci}^h s_{\alpha;t;i},$$

$$\tilde{\Omega}_{\alpha;tt'}^h \equiv \frac{1}{\sqrt{C}} \sum_{c=1}^C q_{\alpha;t;c}^h k_{\alpha;t';c}^h$$

Decoder:

$$\Omega_{\alpha;tt'}^h \equiv \begin{cases} \frac{\exp(\tilde{\Omega}_{\alpha;tt'}^h)}{\sum_{t''=1}^t \exp(\tilde{\Omega}_{\alpha;tt''}^h)} & \text{for } t' \leq t, \\ 0 & \text{for } t' > t, \end{cases}$$

4. Multi-Head Self-Attention Blocks

$$\mathbb{R}_{t;i}^{\text{MHSA}}(s_\alpha; \theta)$$

$$q_{\alpha;t;c}^h \equiv \sum_{i=1}^n Q_{ci}^h s_{\alpha;t;i},$$

$$k_{\alpha;t;c}^h \equiv \sum_{i=1}^n K_{ci}^h s_{\alpha;t;i},$$

$$v_{\alpha;t;c}^h \equiv \sum_{i=1}^n V_{ci}^h s_{\alpha;t;i},$$

$$\tilde{\Omega}_{\alpha;tt'}^h \equiv \frac{1}{\sqrt{C}} \sum_{c=1}^C q_{\alpha;t;c}^h k_{\alpha;t';c}^h$$

Decoder: $\Omega_{\alpha;tt'}^h \equiv \begin{cases} \frac{\exp(\tilde{\Omega}_{\alpha;tt'}^h)}{\sum_{t''=1}^t \exp(\tilde{\Omega}_{\alpha;tt''}^h)} & \text{for } t' \leq t, \\ 0 & \text{for } t' > t, \end{cases}$

$$r_{\alpha;t;i} \equiv \sum_{h=1}^H \sum_{c=1}^C U_{ic}^h \left(\sum_{t'=1}^T \Omega_{\alpha;tt'}^h v_{\alpha;t';c}^h \right)$$

5. Head

If we tie weights between word embedding and word unembedding, then

$$z_{\alpha;t;i}^{(L)} = \mathcal{N}_{\text{rescale}} \sum_{j=1}^n (W^{\text{WE}})^{\top}_{ij} s_{\alpha;t;j}^{(L-1)} = \mathcal{N}_{\text{rescale}} \sum_{j=1}^n W_{ji}^{\text{WE}} s_{\alpha;t;j}^{(L-1)} \quad \text{for } i = 1, \dots, n_{\text{out}},$$

If we don't, then

$$z_{\alpha;t;i}^{(L)} = W_{ij}^{\text{WU}} s_{\alpha;t;j}^{(L-1)}$$