

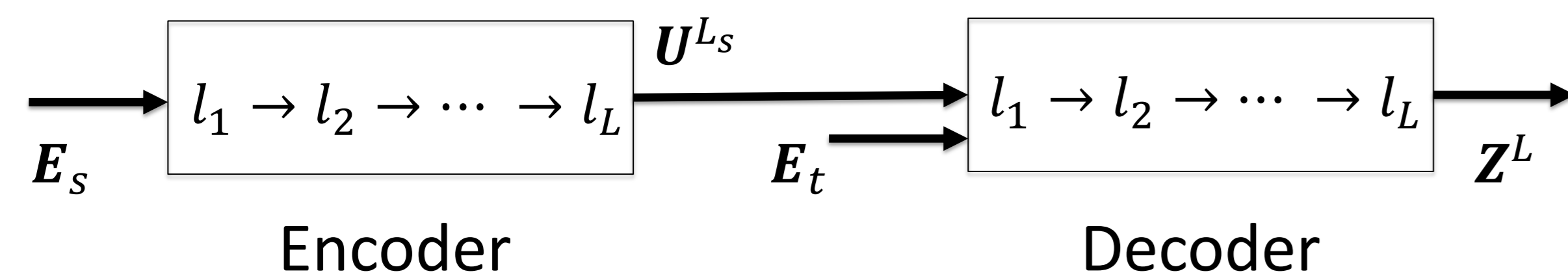
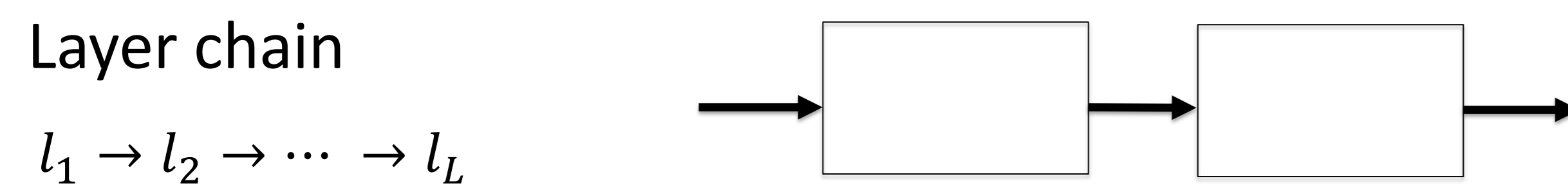
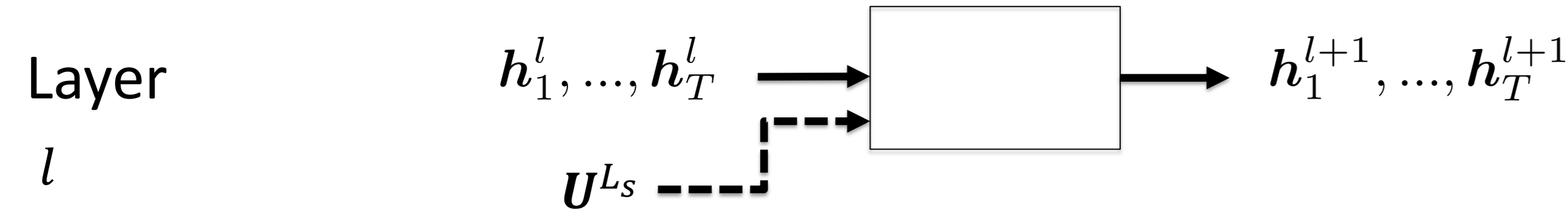
How Much Attention Do You Need? A Granular Analysis of Neural Machine Translation Architectures

Tobias Domhan domhant@amazon.de

Introduction

- **Architecture Definition Language**
- Allows easy experimentation of **Neural Machine Translation (NMT) architecture variations**
- Detailed analysis of NMT architecture combinations
- **Multiple source attention layers and residual feed-forward layer are key**
- **Self-attention more important on the source side**

Architecture Definition Language (ADL)



Layer	Description
<i>pos</i>	Positional embeddings.
<i>norm, bnorm</i>	Layer and batch normalization.
<i>repeat(n, l)</i>	Layer repetition.
<i>res(l)</i>	Residual connection.
<i>birnn, rnn</i>	Recurrent neural network.
<i>cnn</i>	Convolutional neural network.
<i>ff, linear, ffl</i>	Feed-forward and linear layer.
<i>[mh]_dot_src_att(h)</i>	[Multi-head] dot source attention.
<i>[mh]_dot_self_att(h)</i>	[Multi-head] dot self-attention.
<i>res_nd(l)</i>	<i>res(norm → l → dropout)</i>
<i>res_d(l)</i>	<i>res(l → dropout)</i>
<i>act, dropout</i>	Other layers.
<i>parallel, highway</i>	

NMT Architectures defined by ADL

e.g. Transformer [Vaswani et al. 2017]

$$t_{enc} = res_nd(mh_dot_self_att) \rightarrow res_nd(fft)$$

$$t_{dec} = res_nd(mh_dot_self_att) \rightarrow res_nd(mh_dot_src_att) \rightarrow res_nd(fft)$$

$$U^{L_s} = pos \rightarrow repeat(n, t_{enc}) \rightarrow norm$$

$$Z^L = pos \rightarrow repeat(n, t_{dec}) \rightarrow norm$$

Usage in Sockeye, our NMT toolkit:

```
> sockeye-train -s train.de -t train.en
--custom-seq-encoder repeat(6,birnn)
--custom-seq-decoder pos->res(norm->mh_dot_self_att)->res(mh_dot_src_att)
```

<https://github.com/aws-labs/sockeye/tree/acl18>

Encoder/Decoder Combinations

- Baseline: 6 layer Transformer, 512 hidden units
- Data sets: IWSLT'16, WMT'17
- Metrics: BLEU (and METEOR in the paper)
- 3 runs, reporting mean and standard deviation

$$t_{enc} = res_nd() \rightarrow res_nd(fft)$$

$$t_{dec} = res_nd() \rightarrow res_nd(mh_dot_src_att) \rightarrow res_nd(fft)$$

Encoder	Decoder	IWSLT EN→DE	WMT'17 EN→DE	WMT'17 LV→EN
self-att	self-att	25.4 ± 0.2	27.6 ± 0.0	18.3 ± 0.0
self-att	RNN	25.1 ± 0.1	27.4 ± 0.1	18.4 ± 0.2
self-att	CNN	25.4 ± 0.4	27.6 ± 0.2	18.0 ± 0.3
RNN	self-att	25.8 ± 0.1	27.2 ± 0.1	17.8 ± 0.1
CNN	self-att	25.7 ± 0.1	26.6 ± 0.3	16.8 ± 0.4
RNN	RNN	25.1 ± 0.1	26.7 ± 0.1	17.8 ± 0.1
CNN	CNN	25.3 ± 0.3	26.9 ± 0.1	16.4 ± 0.2
self-att	<i>combined</i>	25.1 ± 0.2	27.6 ± 0.2	18.3 ± 0.1
self-att	<i>none</i>	23.7 ± 0.2	25.3 ± 0.2	15.9 ± 0.1

Network Structure Experiments

RNN to Transformer

$$U^{L_s} = dropout \rightarrow res_d(birnn) \rightarrow repeat(5, res_d(rnn))$$

$$Z^L = dropout \rightarrow repeat(6, res_d(rnn)) \rightarrow res_d(dot_src_att) \rightarrow res_d(fft)$$

Model	IWSLT EN→DE	WMT'17 EN→DE	WMT'17 LV→EN
Transformer	25.4 ± 0.1	27.6 ± 0.0	18.5 ± 0.0
RNMT	23.2 ± 0.2	25.5 ± 0.2	-
- input feeding	23.1 ± 0.2	24.6 ± 0.1	-
RNN	22.8 ± 0.2	23.8 ± 0.1	15.2 ± 0.1
+ mh	23.7 ± 0.4	24.4 ± 0.1	16.0 ± 0.1
+ pos	23.9 ± 0.2	24.1 ± 0.1	15.6 ± 0.1
+ norm	23.7 ± 0.1	24.0 ± 0.2	15.2 ± 0.1
+ multi-att-1h	24.5 ± 0.0	25.2 ± 0.1	16.6 ± 0.2
/ multi-att	24.4 ± 0.3	25.5 ± 0.0	17.0 ± 0.2
+ ff	25.1 ± 0.1	26.7 ± 0.1	17.8 ± 0.1

$$U^{L_s} = pos \rightarrow res_nd(birnn) \rightarrow res_nd(fft) \rightarrow repeat(5, res_nd(rnn)) \rightarrow res_nd(fft) \rightarrow norm$$

$$Z^L = pos \rightarrow repeat(6, res_nd(rnn)) \rightarrow res_nd(mh_dot_src_att) \rightarrow res_nd(fft) \rightarrow norm$$

CNN to Transformer

$$U^{L_s} = pos \rightarrow repeat(6, res_d(cnn))$$

$$Z^L = pos \rightarrow repeat(6, res_d(cnn)) \rightarrow res_d(dot_src_att)$$

Model	IWSLT EN→DE	WMT'17 EN→DE	WMT'17 LV→EN
Transformer	25.4 ± 0.1	27.6 ± 0.0	18.5 ± 0.0
CNN GLU	24.3 ± 0.4	25.0 ± 0.3	16.0 ± 0.5
+ norm	24.1 ± 0.1	-	16.1 ± 0.2
+ mh	24.2 ± 0.2	25.4 ± 0.1	16.1 ± 0.1
+ ff	25.3 ± 0.1	26.8 ± 0.1	16.4 ± 0.2
CNN ReLU	23.6 ± 0.3	23.9 ± 0.1	15.4 ± 0.1
+ norm	24.3 ± 0.1	24.3 ± 0.2	16.0 ± 0.2
+ mh	24.2 ± 0.2	24.9 ± 0.1	16.1 ± 0.1
+ ff	25.3 ± 0.3	26.9 ± 0.1	16.4 ± 0.2

$$U^{L_s} = pos \rightarrow repeat(6, res_nd(cnn)) \rightarrow res_nd(fft) \rightarrow norm$$

$$Z^L = pos \rightarrow repeat(6, res_nd(cnn)) \rightarrow res_nd(mh_dot_src_att) \rightarrow res_nd(fft) \rightarrow norm$$