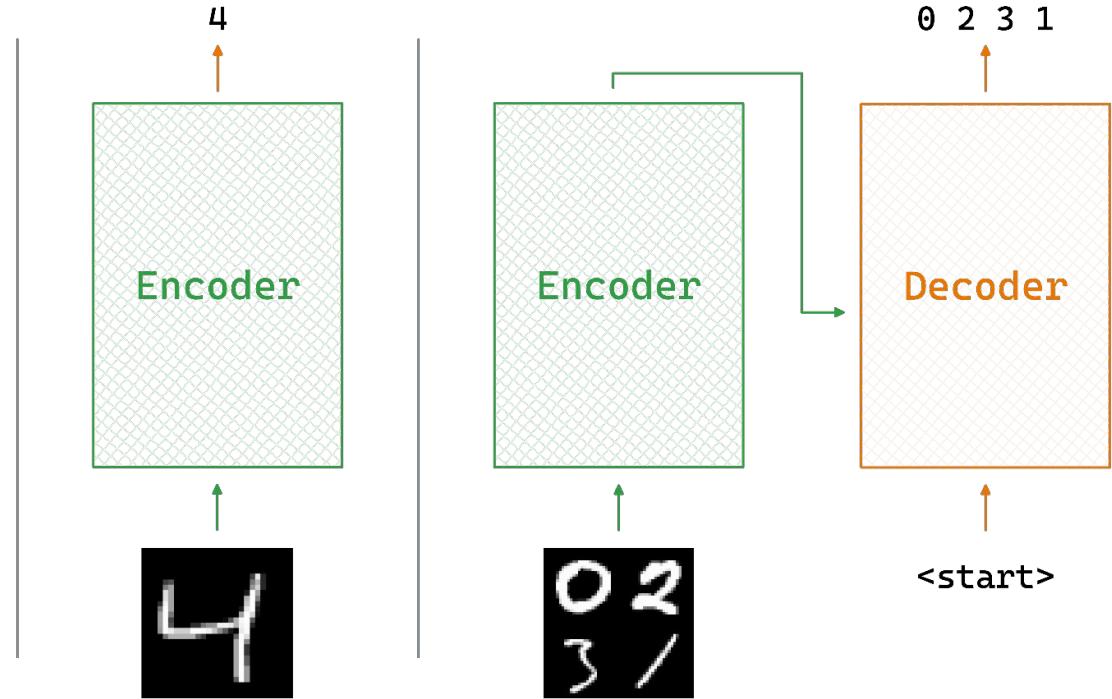


Some Code Is All You Need

A bit of a Q&A

Task

0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9



Paper

Provided proper attribution is provided, Google hereby grants permission to reproduce the tables and figures in this paper solely for use in journalistic or scholarly works.

Attention Is All You Need

Ashish Vaswani*
 Google Brain
 avaswani@google.com

Noam Shazeer*
 Google Brain
 noam@google.com

Niki Parmar*
 Google Research
 nixp@google.com

Jakob Uszkoreit*
 Google Research
 usz@google.com

Llion Jones*
 Google Research
 llion@google.com

Aidan N. Gomez[†]
 University of Toronto
 aidan@cs.toronto.edu

Lukasz Kaiser^{*}
 Google Brain
 lukasz.kaiser@google.com

Illia Polosukhin[‡]
 illia.polosukhin@gmail.com

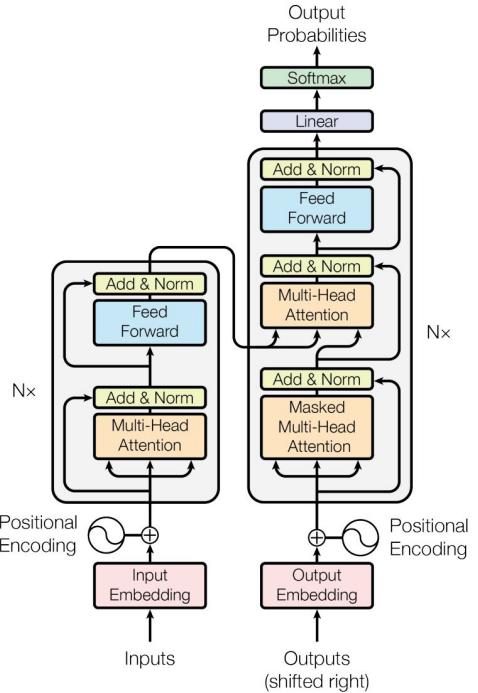
Abstract

The dominant sequence-to-sequence translation models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network structure, the Transformer, based solely on attention mechanisms, without recurrent layers or convolutional layers entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the previous best using neural ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training time required by the best prior work in literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

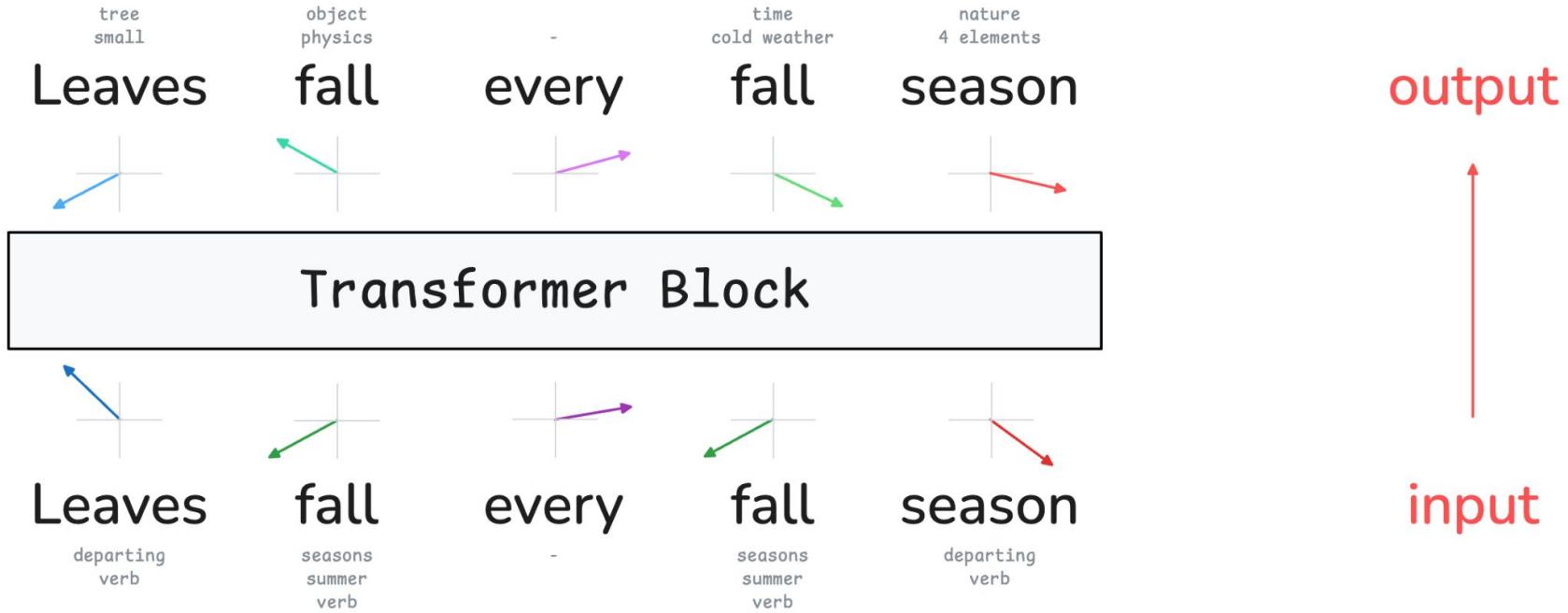
*Equal contribution. [†]Using older codebase. [‡]Work performed while at Google Brain.
[†]Work performed while at Google Research.

31st Conference on Neural Information Processing Systems (NIPS 2017), Long Beach, CA, USA.

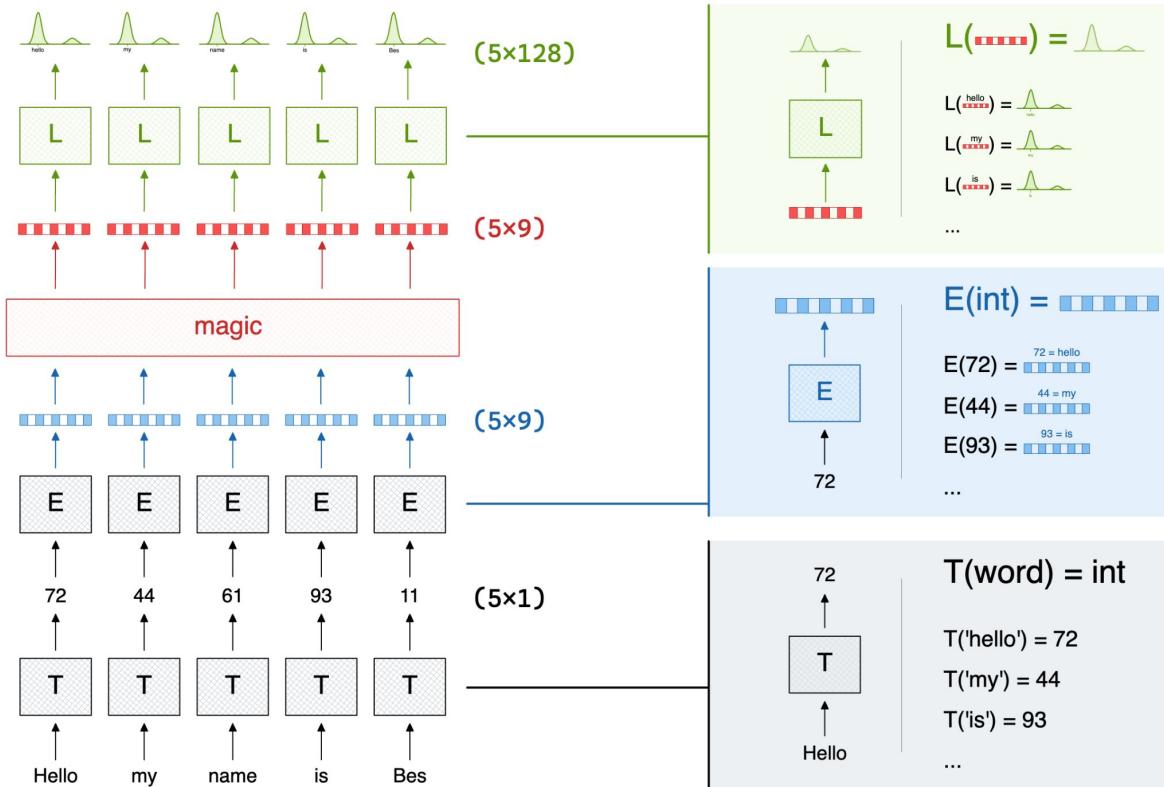
Vaswani et al. 2017



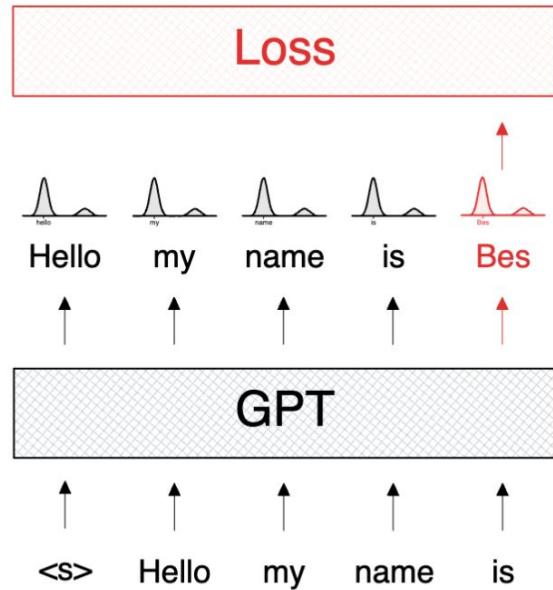
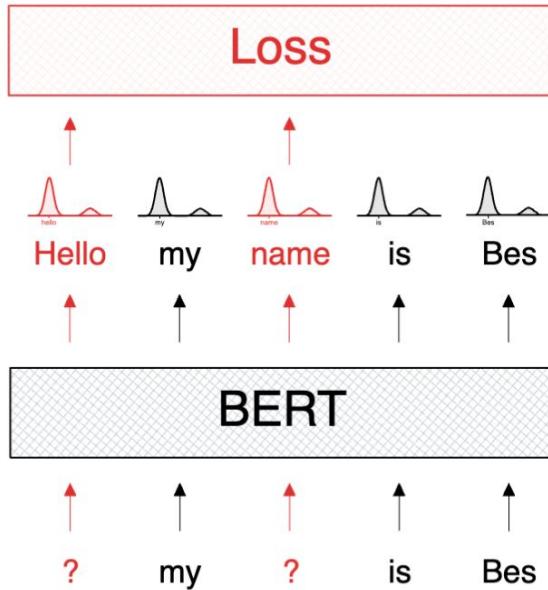
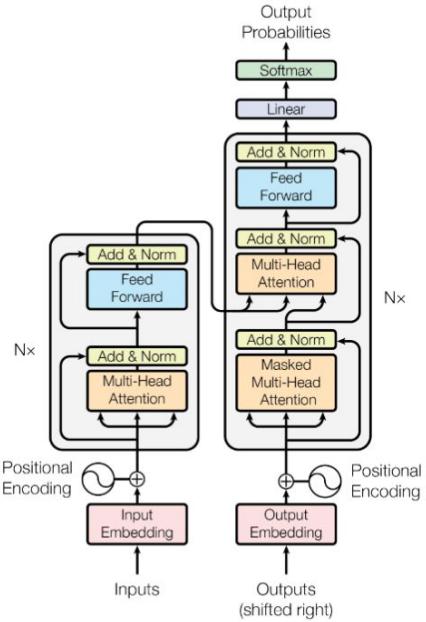
Why “transformer”?



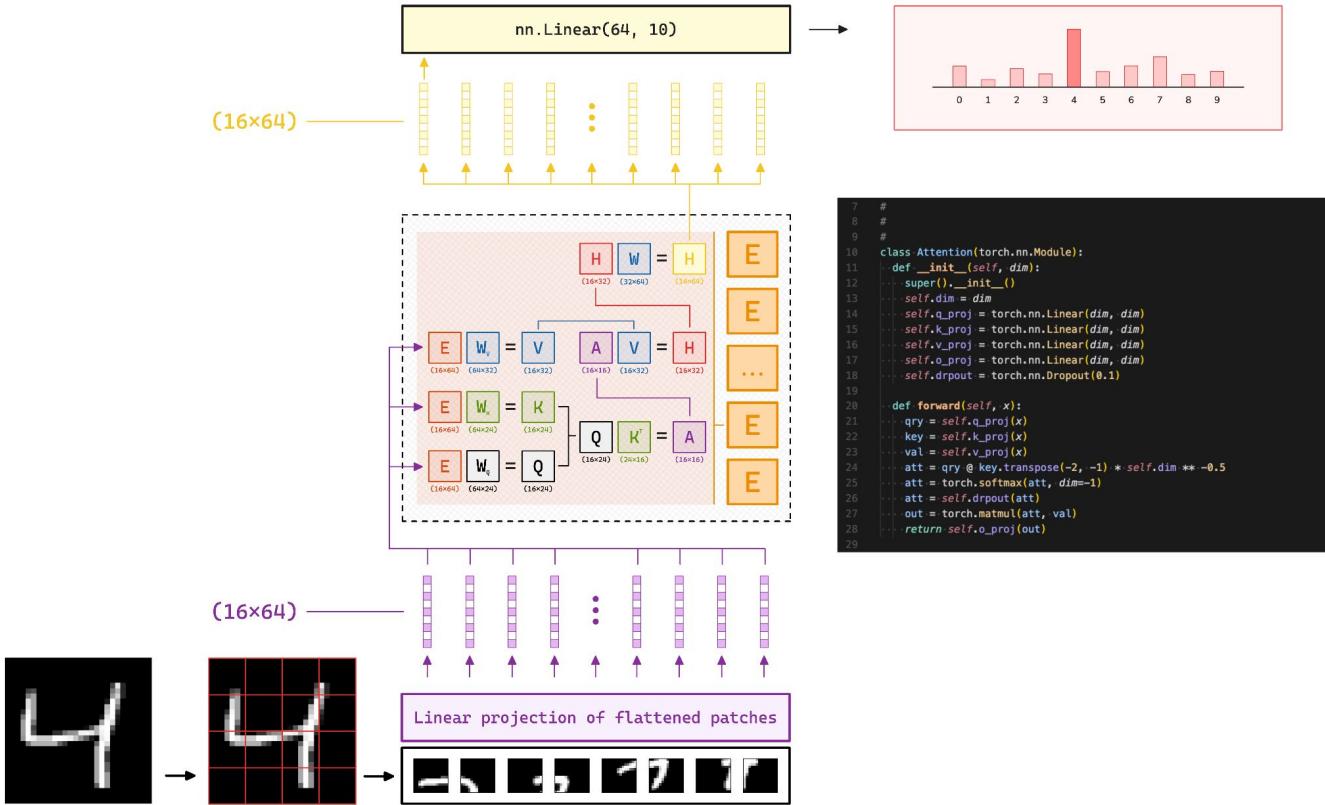
word2vec → transformer



Two Flavours



Encoder



```

7 # 
8 # 
9 # 
10 class Attention(torch.nn.Module):
11     def __init__(self, dim):
12         super().__init__()
13         self.dim = dim
14         self.q_proj = torch.nn.Linear(dim, dim)
15         self.k_proj = torch.nn.Linear(dim, dim)
16         self.v_proj = torch.nn.Linear(dim, dim)
17         self.o_proj = torch.nn.Linear(dim, dim)
18         self.dropout = torch.nn.Dropout(0.1)
19
20     def forward(self, x):
21         qry = self.q_proj(x)
22         key = self.k_proj(x)
23         val = self.v_proj(x)
24         att = qry @ key.transpose(-2, -1) * self.dim ** -0.5
25         att = torch.softmax(att, dim=-1)
26         att = self.dropout(att)
27         out = torch.matmul(att, val)
28         return self.o_proj(out)
29

```

Code

```

72 #
73 #
74 #
75 class LookerTrns(torch.nn.Module):
76     def __init__(self):
77         super().__init__()
78         self.cls = torch.nn.Parameter(torch.randn(1, 1, 128))
79         self.emb = torch.nn.Linear(196, 128)
80         self.pos = torch.nn.Embedding(17, 128)
81         self.register_buffer('rng', torch.arange(17))
82         self.enc = torch.nn.ModuleList([EncoderLayer(128) for _ in range(6)])
83         self.fin = torch.nn.Sequential(
84             torch.nn.LayerNorm(128),
85             torch.nn.Linear(128, 10)
86         )
87
88     def forward(self, x):
89         B = x.shape[0] ..... # [B, 16, 196]
90         pch = self.emb(x) ..... # [B, 16, 128]
91         cls = self.cls.expand(B, -1, -1) .. # [B, 1, 128]
92         hdn = torch.cat([cls, pch], dim=1) .. # [B, 17, 128]
93         hdn = hdn + self.pos(self.rng) .. # [B, 17, 128]
94         for enc in self.enc: hdn = enc(hdn) # [B, 17, 128]
95         out = hdn[:, 0, :] ..... # [B, 128]
96         return self.fin(out) ..... # [B, 10]
97

```

```

50 #
51 #
52 #
53 #
54 class EncoderLayer(torch.nn.Module):
55     def __init__(self, dim):
56         super().__init__()
57         self.att = Attention(dim)
58         self.ffn = FFN(dim)
59         self.ini = torch.nn.LayerNorm(dim)
60         self.fin = torch.nn.LayerNorm(dim)
61
62     def forward(self, src):
63         out = self.att(src)
64         src = src + out
65         src = self.ini(src)
66         out = self.ffn(src)
67         src = src + out
68         src = self.fin(src)
69         return src
70

```

```

31 #
32 #
33 #
34 class FFN(torch.nn.Module):
35     def __init__(self, dim):
36         super().__init__()
37         self.one = torch.nn.Linear(dim, dim)
38         self.drp = torch.nn.Dropout(0.1)
39         self.rlu = torch.nn.ReLU(inplace=True)
40         self.two = torch.nn.Linear(dim, dim)
41
42     def forward(self, x):
43         x = self.one(x)
44         x = self.rlu(x)
45         x = self.drp(x)
46         x = self.two(x)
47         return x
48

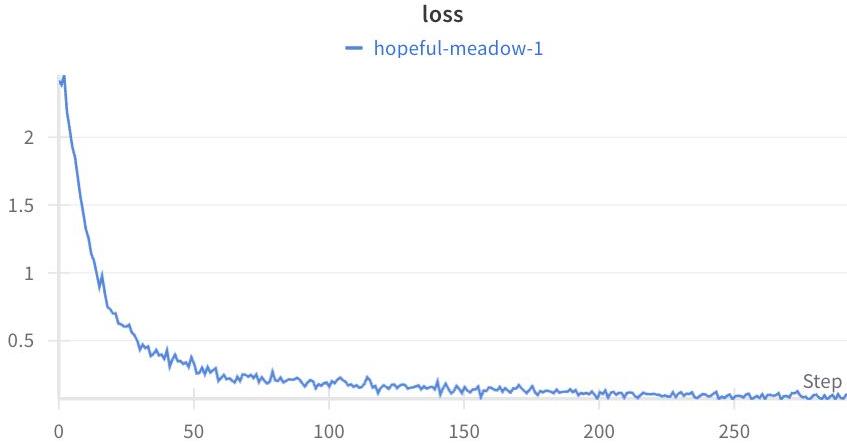
```

```

7 #
8 #
9 #
10 class Attention(torch.nn.Module):
11     def __init__(self, dim):
12         super().__init__()
13         self.dim = dim
14         self.q_proj = torch.nn.Linear(dim, dim)
15         self.k_proj = torch.nn.Linear(dim, dim)
16         self.v_proj = torch.nn.Linear(dim, dim)
17         self.o_proj = torch.nn.Linear(dim, dim)
18         self.dropout = torch.nn.Dropout(0.1)
19
20     def forward(self, x):
21         qry = self.q_proj(x)
22         key = self.k_proj(x)
23         val = self.v_proj(x)
24         att = qry @ key.transpose(-2, -1) * self.dim ** -0.5
25         att = torch.softmax(att, dim=-1)
26         att = self.dropout(att)
27         out = torch.matmul(att, val)
28         return self.o_proj(out)
29

```

It works!



<https://github.com/MLX-FAC/vit>



Thank you!