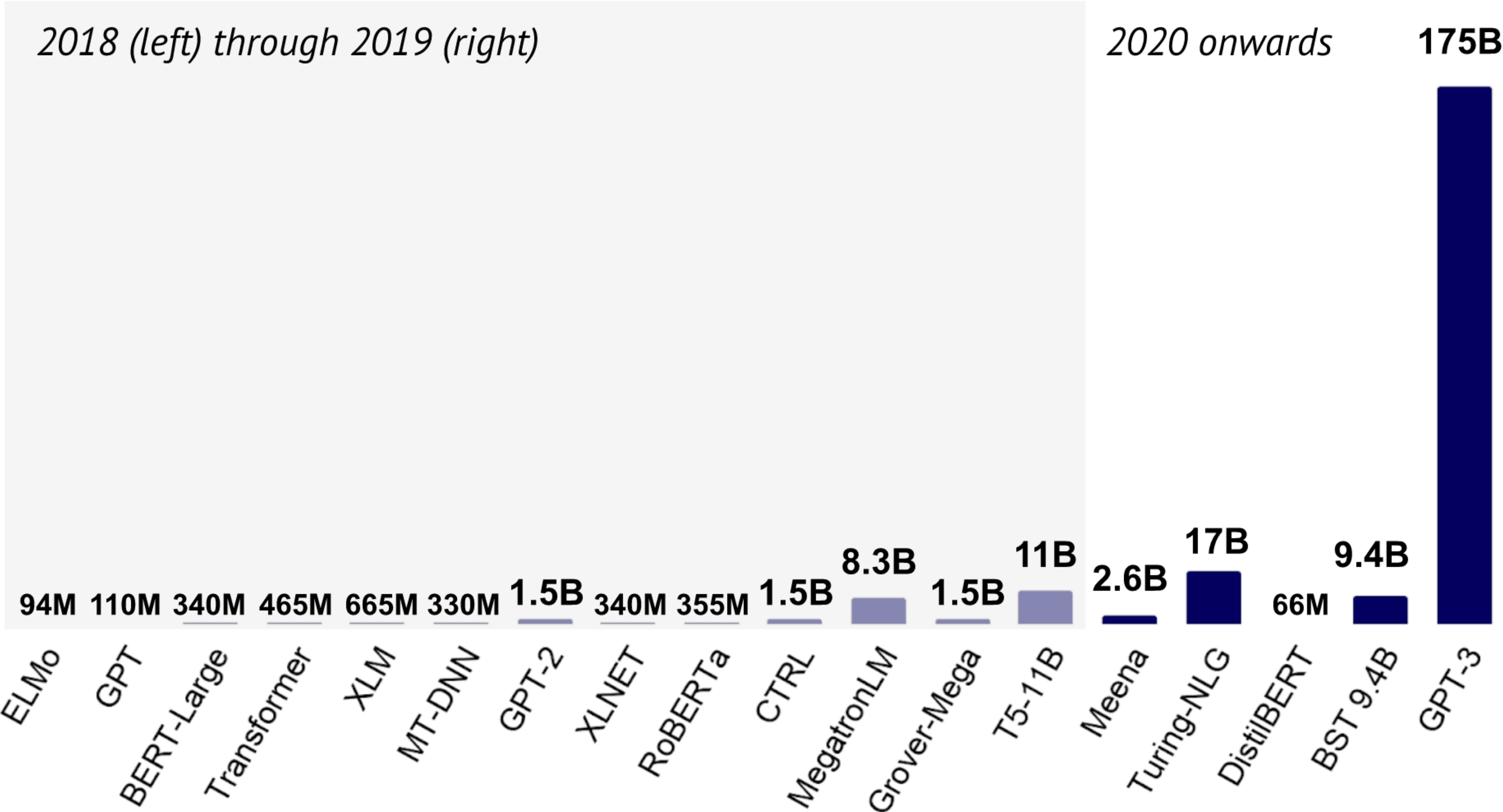


Large Product Key Memory for Pretrained Language Models

**QRAFT | AXE
GUIJIN SON**

Recent NLP Trends

Large Product Key Memory for Pretrained Language Models



Decoupling Parameters & Performance

Large Product Key Memory for Pretrained Language Models

| Model | # Layers | # Params | Inference Speed (batch/sec) |
|-----------------------------|----------|-------------|--------------------------------|
| BERT _{BASE} | 12 | 110M | 79.8 |
| BERT _{BASE} +PKM | 12 | 506M | 61.4 |
| BERT _{BASE} +ResM | 12 | 515M | 59.3 |
| BERT _{LARGE} | 24 | 340M | 43.1 |
| BERT _{LARGE} +PKM | 24 | 860M | 37.2 |
| BERT _{LARGE} +ResM | 24 | 876M | 36.1 |

Product Key Memory

Large Product Key Memory for Pretrained Language Models

sub-key set 1

| |
|-------|
| c_1 |
| c_2 |
| c_3 |

sub-key set 2

| |
|--------|
| c'_1 |
| c'_2 |
| c'_3 |

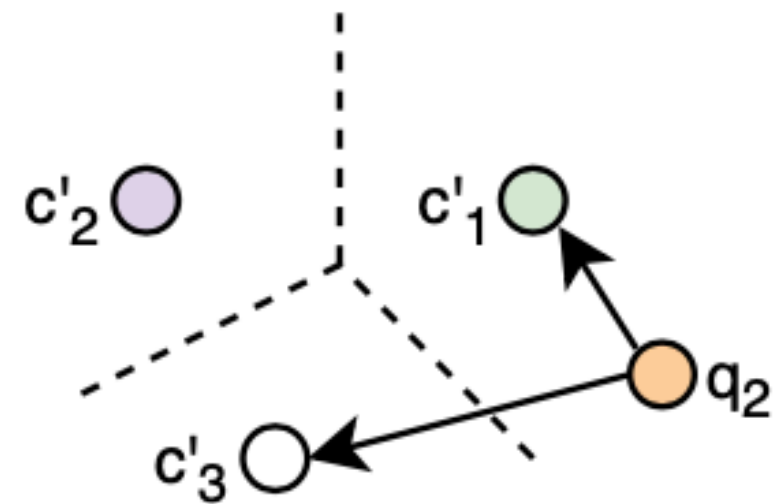
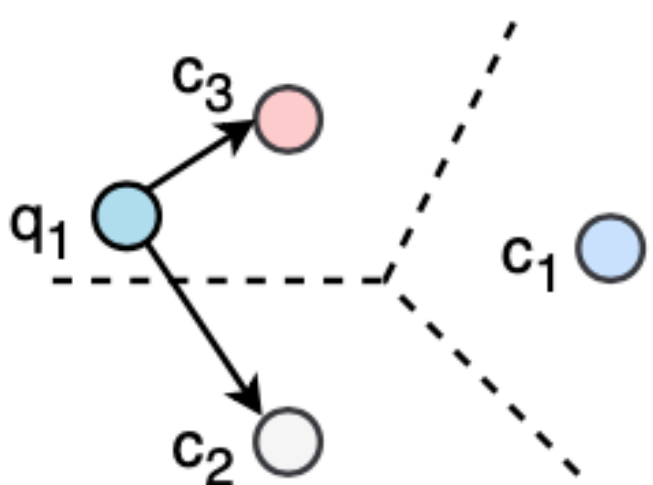
product keys

| | |
|-------|--------|
| c_1 | c'_1 |
| c_1 | c'_2 |
| c_1 | c'_3 |
| c_2 | c'_1 |
| c_2 | c'_2 |
| c_2 | c'_3 |
| c_3 | c'_1 |
| c_3 | c'_2 |
| c_3 | c'_3 |

query



sub-key retrieval



k^2 candidate keys

| | |
|-------|--------|
| c_2 | c'_1 |
| c_2 | c'_3 |
| c_3 | c'_1 |
| c_3 | c'_3 |

key selection

| | |
|-------|--------|
| c_2 | c'_1 |
| c_3 | c'_1 |

k selected keys

Product Key Memory : Pseudo Code

Large Product Key Memory for Pretrained Language Models

```
import torch
import torch.nn as nn
import torch.nn.functional as F

def product_key_memory(query, dim, d_model, d_query, n_head):

    b , t , _ , h = query.shape , n_head

    query = nn.Linear(d_model,d_query)(query)
    query = nn.Dropout(self.MaskedBatchNorm(query))
    query_list = query.chunk(2, dim = -1)

    query = torch.stack(query_list).reshape(2,b,t,h,-1)
    #self.key = nn.parameters(),
    weight = torch.matmul(query_list,self.key,transpose_a = True)
    scores, indices = weight.topk(k=self.topk, dim=-1)
    scores = F.softmax(scores,dim = -1) #normalize

    top_k = sorted(dict((indices,scores)).items(), key=lambda x: x[1],reverse=True)[:10]
    top_k_indice = [indice for indice,score in top_k]

    out = nn.EmbeddingBag(n_key**2,dim,mode='sum')(top_k_indice)

    return nn.Dropout(out)
```

PKM Performance

Large Product Key Memory for Pretrained Language Models

| Model | \widetilde{MU} (4L/8L) (%) | Memory | | WT-2 (ppl) | MLM | |
|---------------------------------------|---------------------------------|-------------------|-------------------|---------------|-----------------|----------------|
| | | KL_u (4L/8L) | KL_w (4L/8L) | | WT-103 (ppl) | PG-19 (ppl) |
| (a) BERT _{BASE} [†] | - | - | - | 3.49 | 3.86 | 6.18 |
| (b) +500k steps | - | - | - | 3.40 | 3.72 | 5.88 |
| (c) +PKM | 2.2/84.1 | 1.62/0.89 | 1.99/1.13 | 3.26 | 3.39 | 5.53 |
| (d) +ResM | 75.0/81.0 | 1.50/0.71 | 1.80/0.92 | 3.26 | 3.36 | 5.45 |
| (e) +Init +PKM | 97.4/95.7 | 0.53/0.69 | 0.68/0.88 | 3.14 | 3.26 | 5.22 |
| (f) +Init +ResM | 98.2/97.3 | 0.45/0.46 | 0.58/0.60 | 3.10 | 3.20 | 5.14 |

Catastrophic Drift

Large Product Key Memory for Pretrained Language Models

| Model | QA | | GLUE | | | | Avg - |
|--|----------------------|----------------------|--------------|---------------|----------------|----------------|-------------|
| | SQuAD 1.1 (EM/F1) | MNLI-(m/mm) (Acc) | QQP (Acc) | QNLI (Acc) | SST-2 (Acc) | CoLA (Matt) | |
| (a) BERT _{BASE} [†] | 82.7/89.8 | 84.3/84.5 | 91.0 | 89.3 | 92.8 | 60.8 | 83.8 |
| (b) +500k steps | 83.3/90.1 | 84.8/84.9 | 91.2 | 89.2 | 92.4 | 61.4 | 84.0 |
| (c) +PKM | 81.9/89.1 | 84.4/85.0 | 91.1 | 89.0 | 93.6 | 59.7 | 83.8 |
| (d) +ResM | 81.5/89.4 | 84.6/84.8 | 91.0 | 88.2 | 93.2 | 62.8 | 84.1 |
| (e) +Init +PKM | 83.8/90.6 | 85.8/85.6 | 91.2 | 90.0 | 93.6 | 63.6 | 85.0 |
| (f) +Init +ResM | 83.9/90.8 | 86.0/85.8 | 91.4 | 90.4 | 94.0 | 64.1 | 85.3 |
| (g) BERT _{BASE} [*] | 81.1/88.5 | 83.9/84.4 | 91.0 | 88.4 | 92.9 | 59.8 | 83.4 |
| (h) BERT _{LARGE} [*] | 83.3/90.6 | 86.2/86.1 | 91.4 | 90.4 | 93.8 | 64.1 | 85.3 |

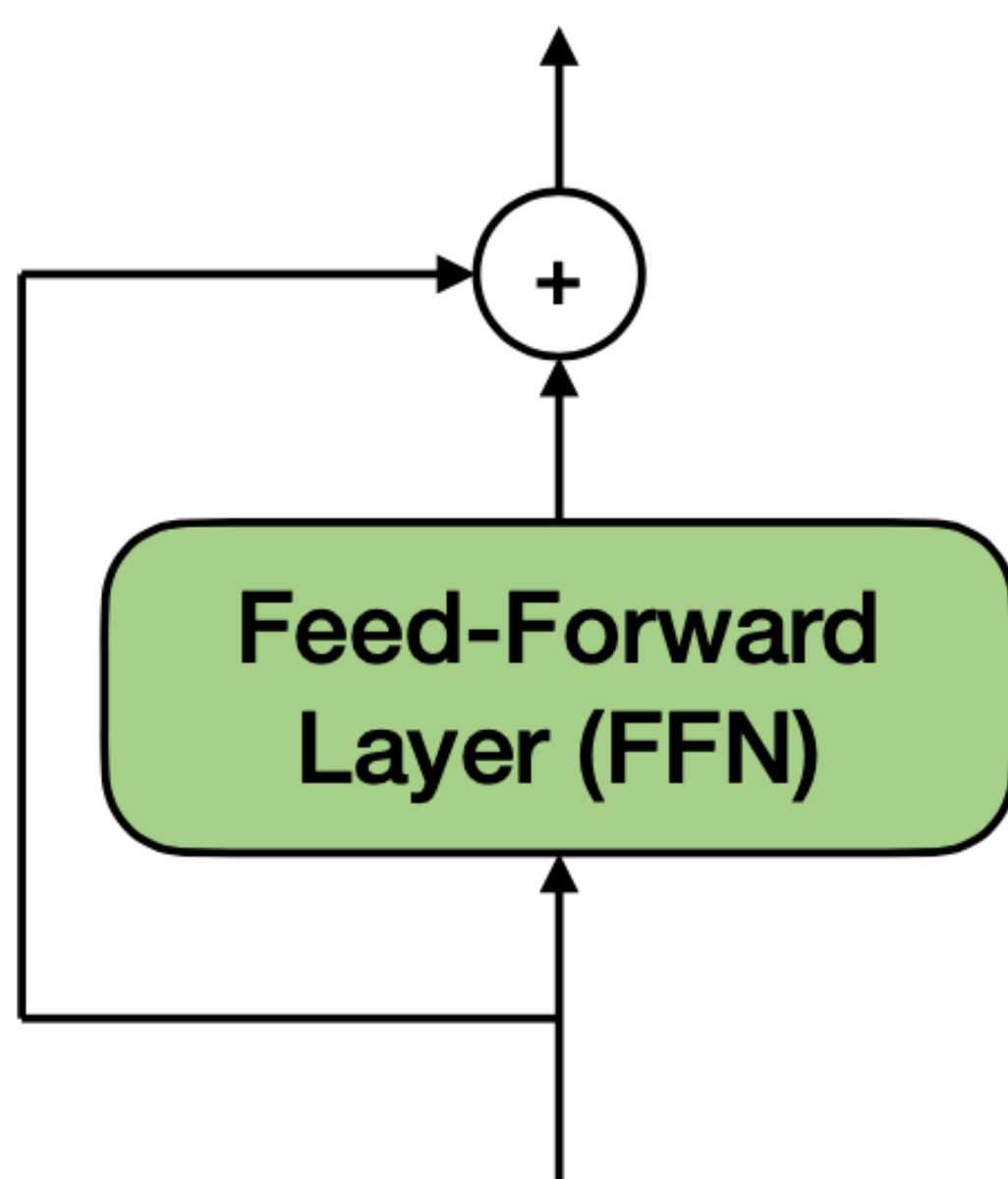
Solutions

Large Product Key Memory for Pretrained Language Models

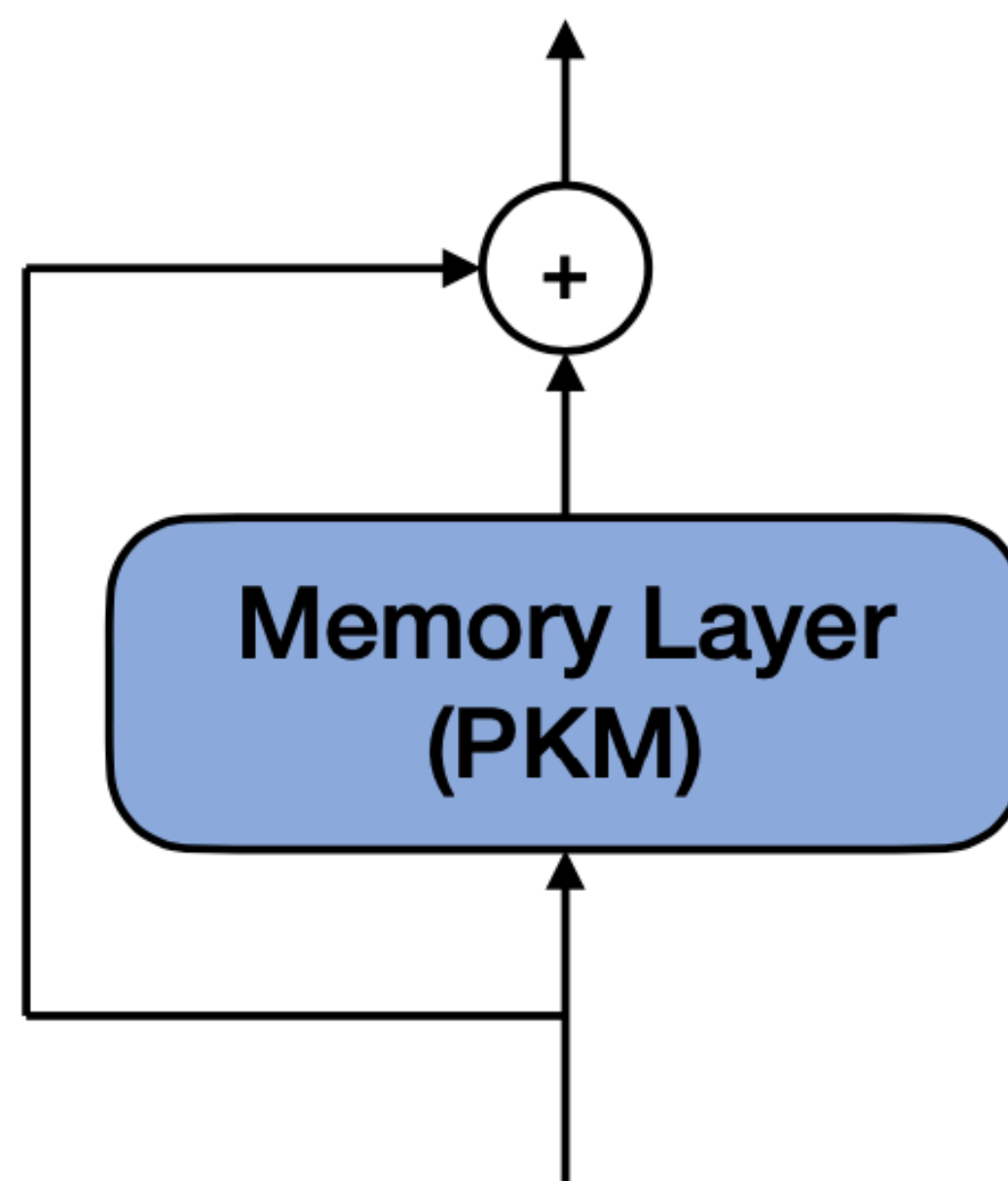
- **Initialization from Pretrained Weights**
- **Residual Memory Layer**

Model Architecture

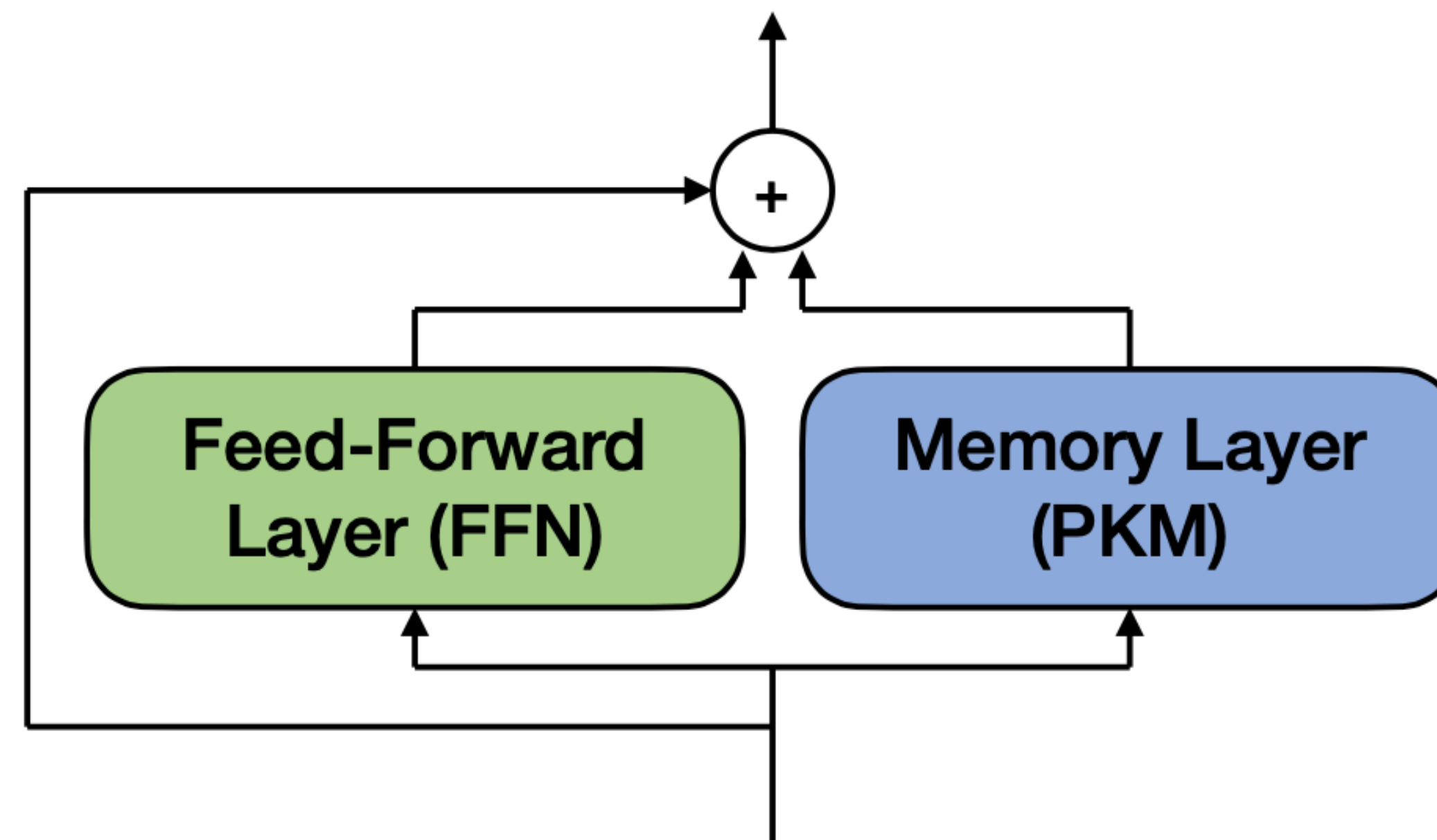
Large Product Key Memory for Pretrained Language Models



(a) FFN layer



(b) PKM layer



(c) ResM layer

ResM Performance

Large Product Key Memory for Pretrained Language Models

| Model | QA | | GLUE | | | | Avg - |
|--|----------------------|----------------------|--------------|---------------|----------------|----------------|-------------|
| | SQuAD 1.1 (EM/F1) | MNLI-(m/mm) (Acc) | QQP (Acc) | QNLI (Acc) | SST-2 (Acc) | CoLA (Matt) | |
| (a) BERT _{BASE} [†] | 82.7/89.8 | 84.3/84.5 | 91.0 | 89.3 | 92.8 | 60.8 | 83.8 |
| (b) +500k steps | 83.3/90.1 | 84.8/84.9 | 91.2 | 89.2 | 92.4 | 61.4 | 84.0 |
| (c) +PKM | 81.9/89.1 | 84.4/85.0 | 91.1 | 89.0 | 93.6 | 59.7 | 83.8 |
| (d) +ResM | 81.5/89.4 | 84.6/84.8 | 91.0 | 88.2 | 93.2 | 62.8 | 84.1 |
| (e) +Init +PKM | 83.8/90.6 | 85.8/85.6 | 91.2 | 90.0 | 93.6 | 63.6 | 85.0 |
| (f) +Init +ResM | 83.9/90.8 | 86.0/85.8 | 91.4 | 90.4 | 94.0 | 64.1 | 85.3 |
| (g) BERT _{BASE} [*] | 81.1/88.5 | 83.9/84.4 | 91.0 | 88.4 | 92.9 | 59.8 | 83.4 |
| (h) BERT _{LARGE} [*] | 83.3/90.6 | 86.2/86.1 | 91.4 | 90.4 | 93.8 | 64.1 | 85.3 |

Further Implications

Large Product Key Memory for Pretrained Language Models

| Model | MLM | QA | GLUE | |
|------------------------|----------------|----------------------|-----------------|----------------|
| | PG-19 (ppl) | SQuAD 1.1 (EM/F1) | MNLI-m (Acc) | SST-2 (Acc) |
| DistilBERT* | 20.61 | 77.4/85.7 | 82.0 | 91.6 |
| +Init +ResM | 5.75 | 80.4/88.3 | 84.1 | 93.3 |
| BERT _{BASE} * | 11.82 | 81.1/88.5 | 83.9 | 92.9 |