

3.1 Unsupervised pre-training

Given an unsupervised corpus of tokens $\mathcal{U} = \{u_1, \dots, u_n\}$, we use a standard language modeling objective to maximize the following likelihood:

$$L_1(\mathcal{U}) = \sum_i \log P(u_i | u_{i-k}, \dots, u_{i-1}; \Theta)$$

where k is the size of the context window, and the conditional probability P is modeled using a neural network with parameters Θ . These parameters are trained using stochastic gradient descent [51].

In our experiments, we use a multi-layer *Transformer decoder* [34] for the language model, which is a variant of the transformer [62]. This model applies a multi-headed self-attention operation over the input context tokens followed by position-wise feedforward layers to produce an output distribution over target tokens:

$$h_0 = UW_e + W_p$$
$$h_i = \text{transformer_block}(h_{i-1}) \forall i \in [1, n]$$
$$P(u) = \text{softmax}(h_n W_e^T)$$

where $U = (u_{-k}, \dots, u_{-1})$ is the context vector of tokens, n is the number of layers, W_e is the token embedding matrix, and W_p is the position embedding matrix.

3.2 Supervised fine-tuning

After training the model with the objective in Eq. 1, we adapt the parameters to the supervised target task. We assume a labeled dataset \mathcal{C} , where each instance consists of a sequence of input tokens, x^1, \dots, x^m , along with a label y . The inputs are passed through our pre-trained model to obtain the final transformer block's activation h_i^m , which is then fed into an added linear output layer with parameters W_y to predict y :

$$P(y|x^1, \dots, x^m) = \text{softmax}(h_i^m W_y).$$

This gives us the following objective to maximize:

$$L_2(\mathcal{C}) = \sum_{(x,y)} \log P(y|x^1, \dots, x^m).$$

We additionally found that including language modeling as an auxiliary objective to the fine-tuning helped learning by (a) improving generalization of the supervised model, and (b) accelerating convergence. This is in line with prior work [50, 43], who also observed improved performance with such an auxiliary objective. Specifically, we optimize the following objective (with weight λ):

$$L_3(\mathcal{C}) = L_2(\mathcal{C}) + \lambda * L_1(\mathcal{C})$$

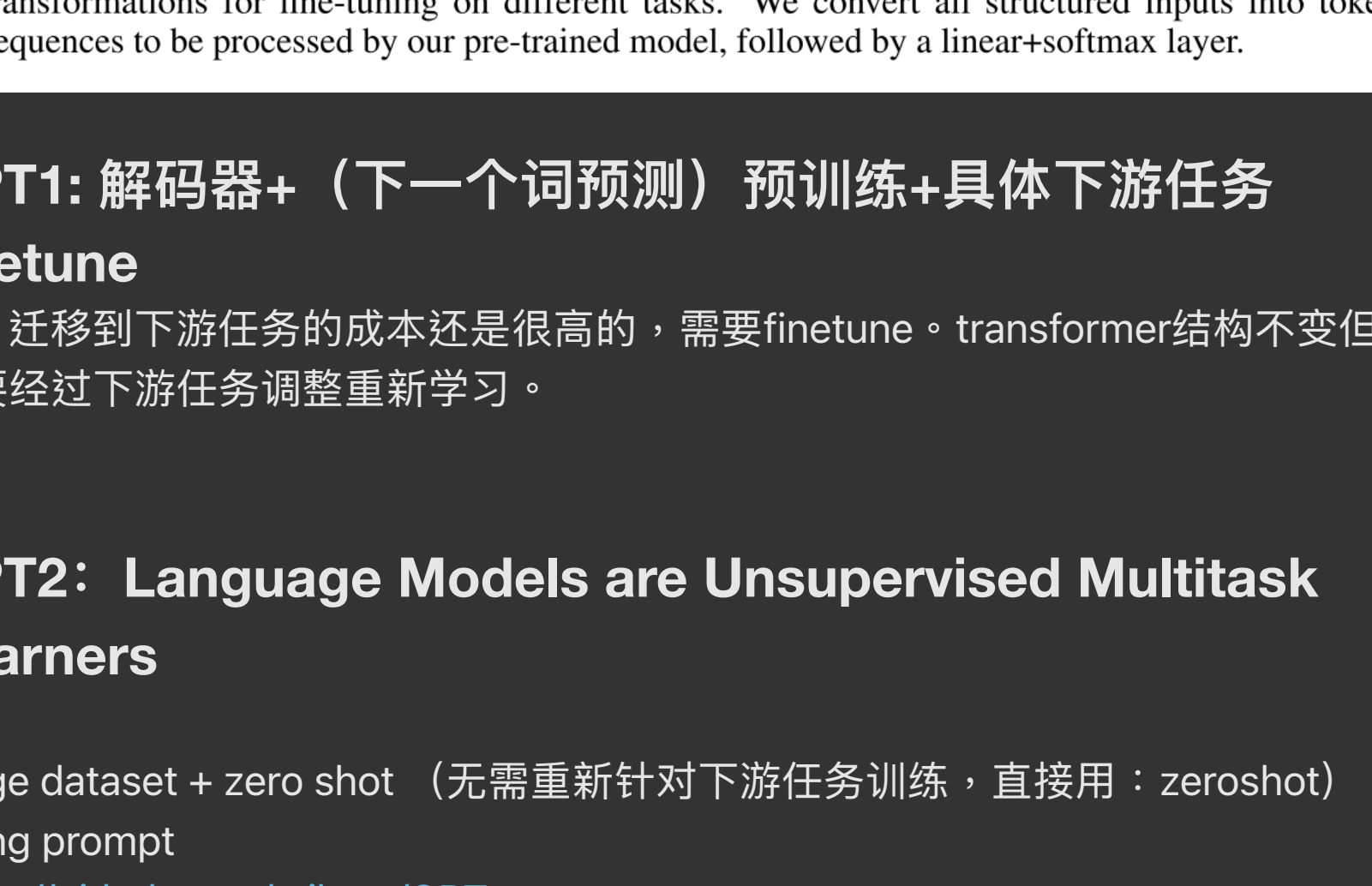


Figure 1: (left) Transformer architecture and training objectives used in this work. (right) Input transformations for fine-tuning on different tasks. We convert all structured inputs into token sequences to be processed by our pre-trained model, followed by a linear+softmax layer.

GPT1: 解码器+（下一个词预测）预训练+具体下游任务 finetune

迁移到下游任务的成本还是很高的，需要finetune。transformer结构不变但参数要经过下游任务调整重新学习。

GPT2: Language Models are Unsupervised Multitask Learners

Large dataset + zero shot（无需重新针对下游任务训练，直接用：zeroshot）

Using prompt

<https://github.com/milomaj/gpt>

GPT2改进之处：

- 去掉下游任务没见过的分隔符之类的特殊符号。引入prompt方法：
 - 翻译任务prompt: (translate to france, English text, French text)
 - 阅读理解任务prompt: (answer the question, document, question, answer)
- 模型初始改变（scale the weights of residual layers at initialization by a factor of $1/\sqrt{N}$ where N is the number of residual layers）
- Add sparse transformer的改进
- pre normalization，reversible tokenization
- 训练数据集构建：
 - 没用Common Crawl 信噪比很低。
 - 采用Reddit，kamma（3个以上的问答数据集）。4500w个link，800w文本，40GB文字。
 - 4个大小不同的模型

GPT3: Language Models are Few-Shot Learners

GPT-3, an autoregressive language model with 175 billion parameters, 10x more than any previous non-sparse language model, and test its performance in the few-shot setting. GPT-3 is applied **without any gradient updates or fine-tuning**. GPT-3 can generate samples of news articles which human evaluators have difficulty distinguishing from articles written by humans.

性能提升：

- weights 训练好之后不会再更新了
- in-context learning
- Evaluation of GPT-3
 - Few shot learning每个sample提供10-100个例子。
 - One shot learning
 - Zero shot learning

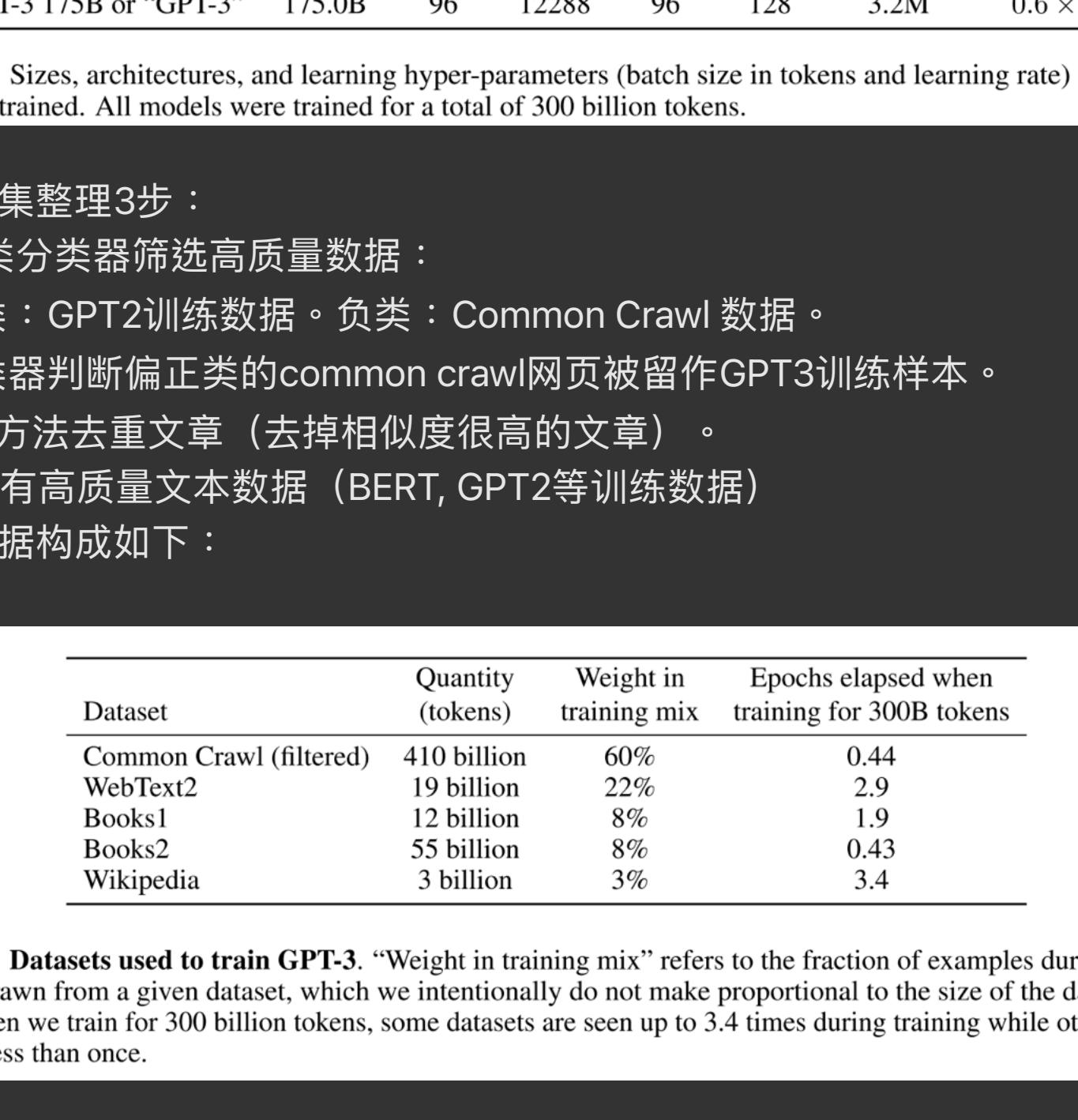


Figure 2.1: Zero-shot, one-shot and few-shot, contrasted with traditional fine-tuning. The panels above show four methods for performing a task with a language model – fine-tuning is the traditional method, whereas zero-, one-, and few-shot, which we study in this work, require the model to perform the task with only forward passes at test time. We typically present the model with a few dozen examples in the few shot setting. Exact phrasings for all task descriptions, examples and prompts can be found in Appendix G.

模型：

Model Name	n_{params}	n_{layers}	d_{model}	n_{heads}	d_{head}	Batch Size	Learning Rate
GPT-3 Small	125M	12	768	12	64	0.5M	6.0×10^{-4}
GPT-3 Medium	350M	24	1024	16	64	0.5M	6.0×10^{-4}
GPT-3 Large	760M	24	1536	16	96	0.5M	2.5×10^{-4}
GPT-3 XL	1.3B	24	2048	24	128	1M	2.0×10^{-4}
GPT-3 2.7B	2.7B	32	2560	32	80	1M	1.6×10^{-4}
GPT-3 6.7B	6.7B	32	4096	32	128	2M	1.2×10^{-4}
GPT-3 13B	13.0B	40	5140	40	128	2M	1.0×10^{-4}
GPT-3 175B or “GPT-3”	175.0B	96	12288	96	128	3.2M	0.6×10^{-4}

Table 2.1: Sizes, architectures, and learning hyper-parameters (batch size in tokens and learning rate) of the models which we trained. All models were trained for a total of 300 billion tokens.

训练数据集整理3步：

- 构建2类分类器筛选高质量数据：
 - 正类：GPT2训练数据。负类：Common Crawl 数据。
 - 分类器判断偏正类的common crawl网页被留作GPT3训练样本。
- 用LSH方法去重文章（去掉相似度很高的文章）。
- 加入已有高质量文本数据（BERT, GPT2等训练数据）

此后的数据构成如下：

Dataset	Quantity (tokens)	Weight in training mix	Epochs elapsed when training for 300B tokens
Common Crawl (filtered)	410 billion	60%	0.44
WebText2	19 billion	22%	2.9
Books1	12 billion	8%	1.9
Books2	55 billion	8%	0.43
Wikipedia	3 billion	3%	3.4

Table 2.2: Datasets used to train GPT-3. “Weight in training mix” refers to the fraction of examples during training that are drawn from a given dataset, which we intentionally do not make proportional to the size of the dataset. As a result, when we train for 300 billion tokens, some datasets are seen up to 3.4 times during training while other datasets are seen less than once.

