

09

实战演练之 生成式对话机器人

- (1) 对话机器人简介
- (2) 代码实战演练，基于BLOOM
- (3) 常见解码参数介绍

对话机器人介绍

对话机器人简介

- 什么是对话机器人

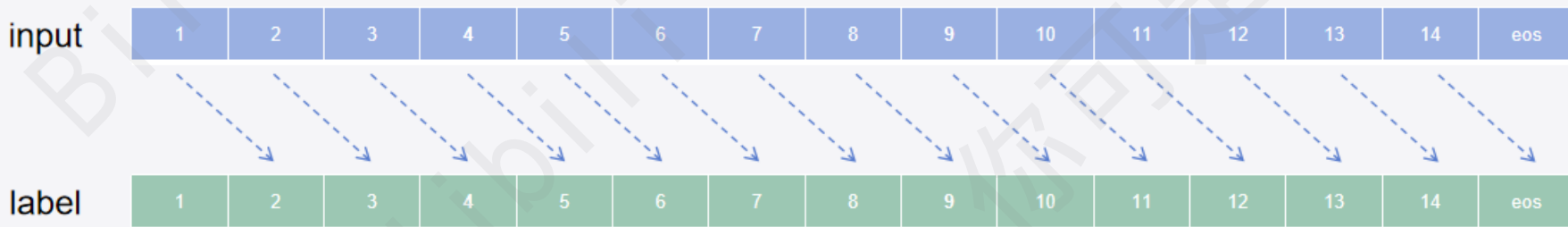
- 对话机器人在本质上是一个用来模拟人类对话或聊天的计算机程序，接收人类的自然语言作为输入，并给出合适的回复
- 按照任务类型划分，对话机器人简单的可以分为闲聊机器人、问答机器人、任务型对话机器人
- 按照答案产生的逻辑划分，对话机器人可以划分为检索式对话机器人和生成式对话机器人
- 本次课程关注的内容为**基于生成式的问答机器人**

对话机器人解决方案

预训练简介

- 预训练任务

- 因果语言模型，自回归模型
 - 将完整序列输入，基于上文的token预测当前token
 - 结束位置要有特殊token, eos_token

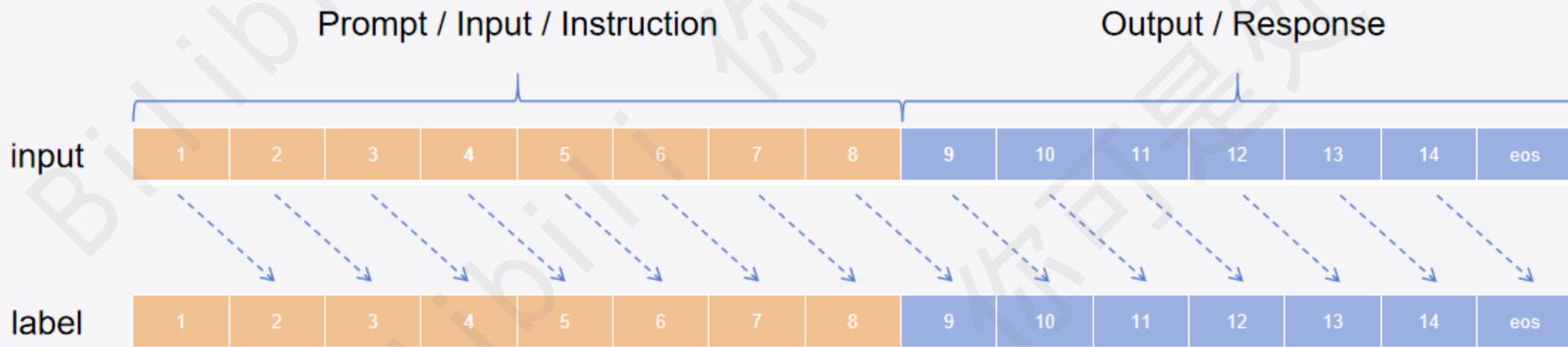


对话机器人解决方案

对话机器人解决方案

- 指令微调

- 指令微调的方式，赋予回答问题的能力
- 多类型的任务共同学习，能够解决不同的任务



对话机器人解决方案

对话机器人解决方案

- 指令微调

- 指令微调的方式，赋予回答问题的能力
 - 训练单轮问答模型，计算Loss时只计算Output部分



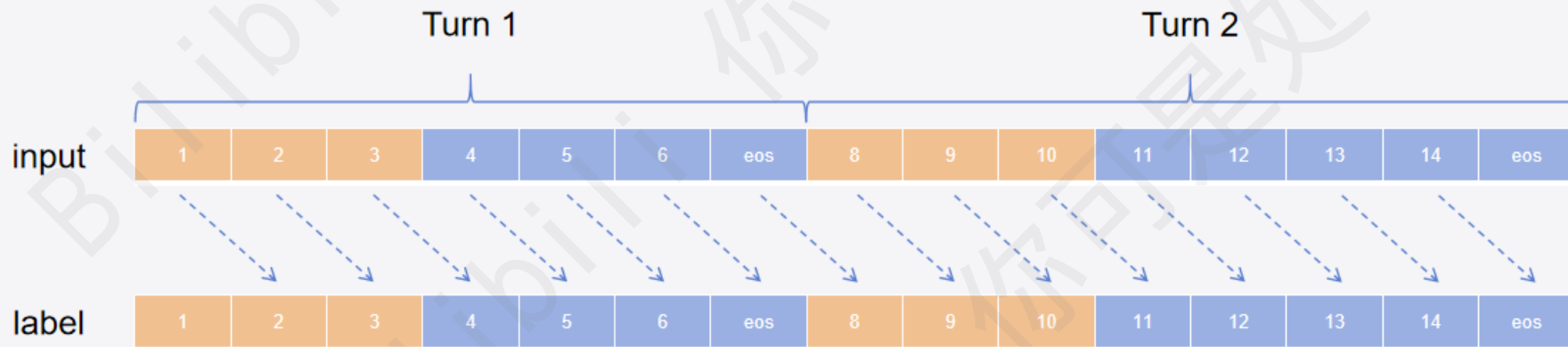
对话机器人解决方案

对话机器人解决方案

- 指令微调

- 指令微调的方式，赋予回答问题的能力

- 多轮如何计算？



对话机器人解决方案

对话机器人解决方案

- 指令微调

- 指令微调的方式，赋予回答问题的能力
 - 方式一 计算最后一轮Output的Loss，效率较低

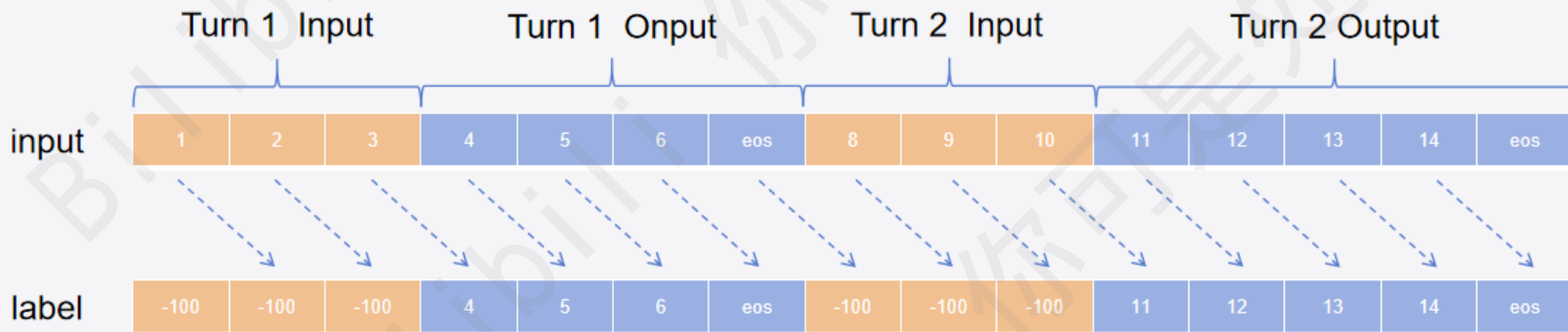


对话机器人解决方案

对话机器人解决方案

- 指令微调

- 指令微调的方式，赋予回答问题的能力
 - 方式二 计算每一轮Output的Loss，效率更高



代码实战演练

代码实战演练（基于Bloom模型）

- 数据集

- <https://huggingface.co/datasets/shibing624/alpaca-zh>
- 指令微调

- 预训练模型

- Langboat/bloom-389m-zh

常见解码参数介绍

常见解码参数介绍

• 常用推理参数

- 长度控制
 - `min/max_new_tokens` 最小/最大生成的长度
 - `min/max_length` 序列整体的最小/最大长度
- 解码策略
 - `do_sample` 是否启用采样的生成方式
 - `num_beams beam_search`的大小
- 采样参数
 - `temperature` 默认1.0，即原始分布，低于1.0会使得分布更尖锐，高于1.0会使得分布更均匀
 - `top_k` 将词概率从大到小排列，将采样限制在前K个词
 - `top_p` 将词概率从大到小排列，将采样限制在前N个词，条件是这N个词的概率超过`top_p`的值
- 惩罚项
 - `repetition_penalty` 重复惩罚项，实现原理是降低已经出现过的token的概率