

天才报大腿

Contents

目录

1 团队介绍

2 赛题理解

3 特征工程

4 模型介绍

5 总结与思考



团队介绍

团队介绍



陶超
华中科技大学
统计学硕士
算法工程师



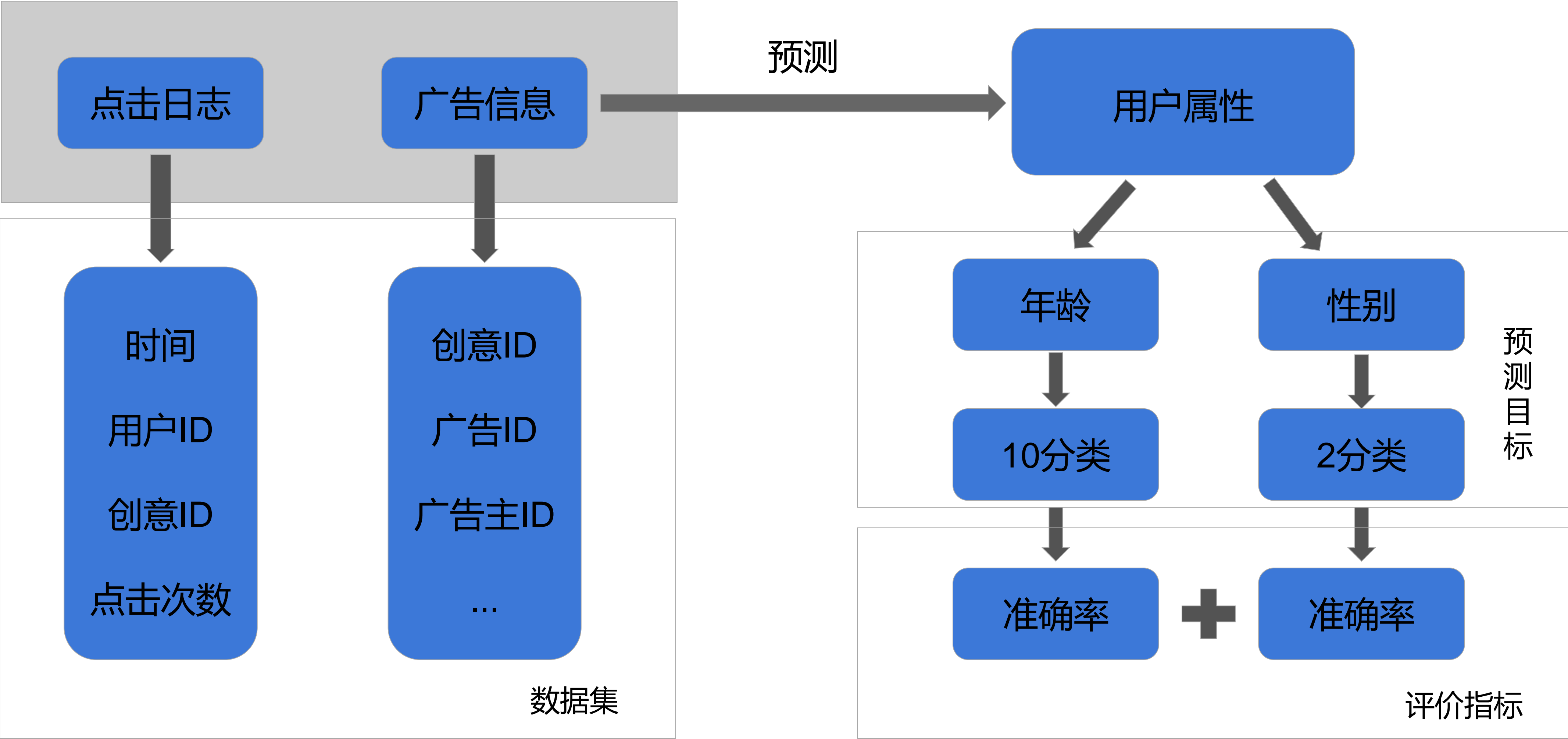
程唯
南京大学
计算机硕士
算法工程师



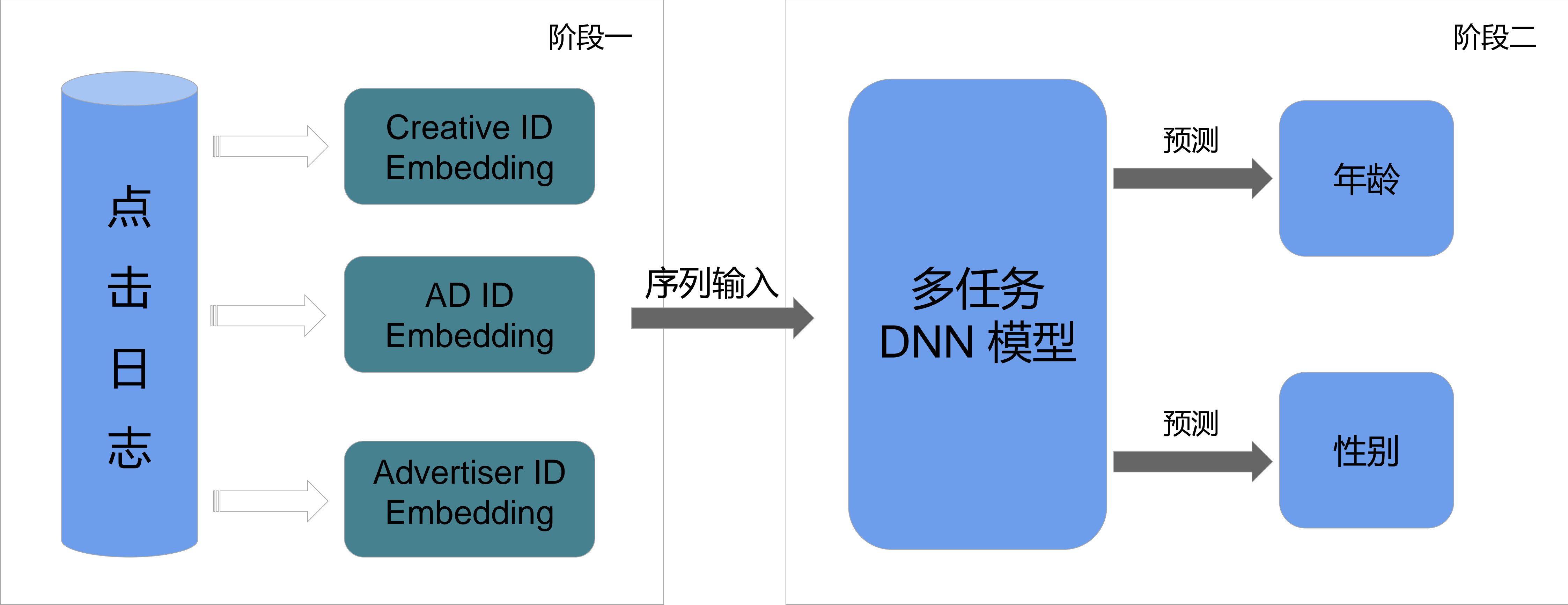
罗欣晨
四川大学
计算机本科
算法工程师



赛题理解



两阶段建模





特征工程

Embedding构建/数据增广

Embedding构建

通过点击日志构造每个用户的浏览序列

User 1

浏览序列

User 2

浏览序列

User 3

浏览序列

User 4

浏览序列

User 5

浏览序列

■ ■ ■

支持对数据在线shuffle

Word2Vec
CBOW

Word2Vec
Skip-Gram

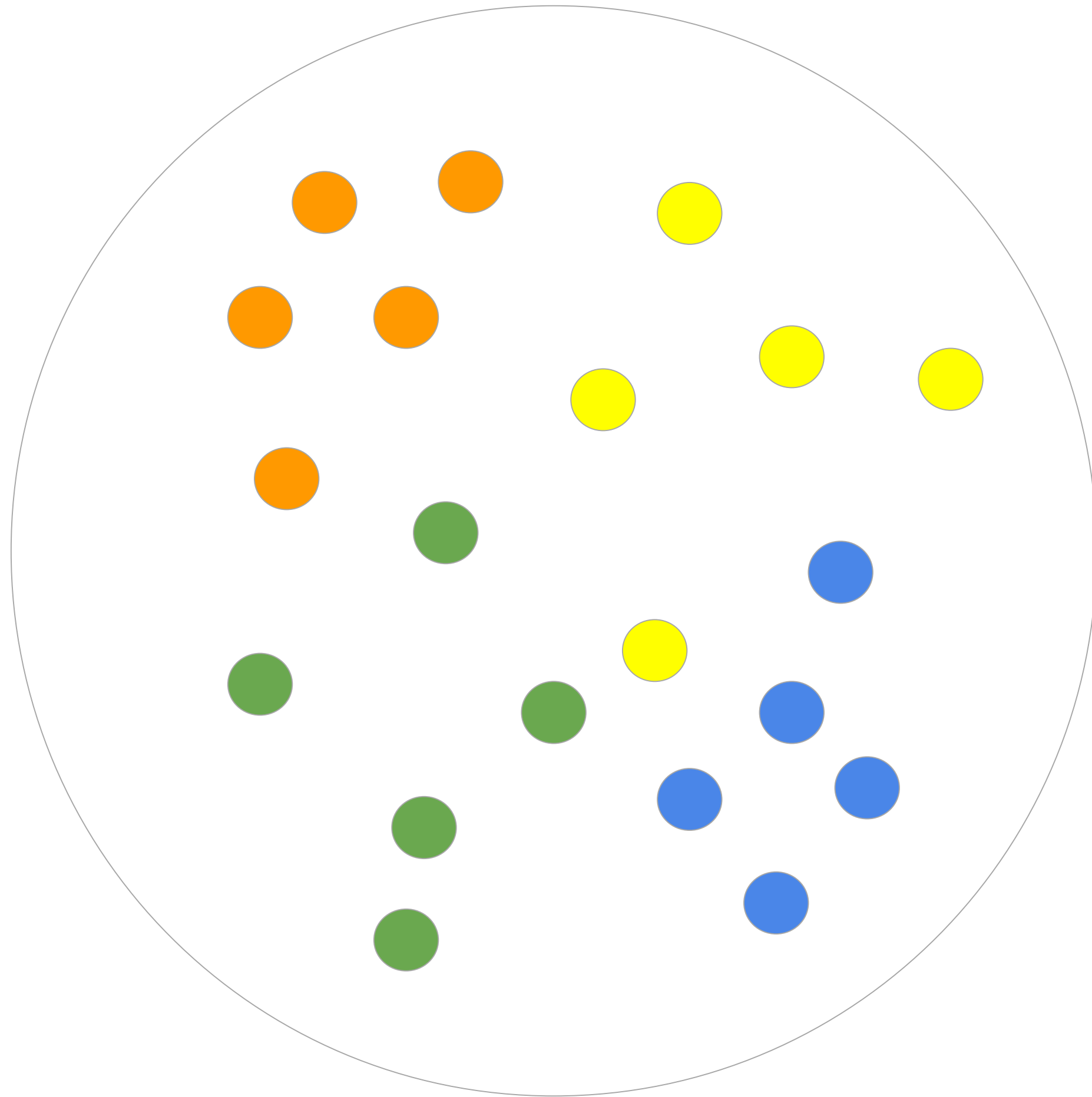
GloVe

CreativeID/
ADID/
AdvertiserID

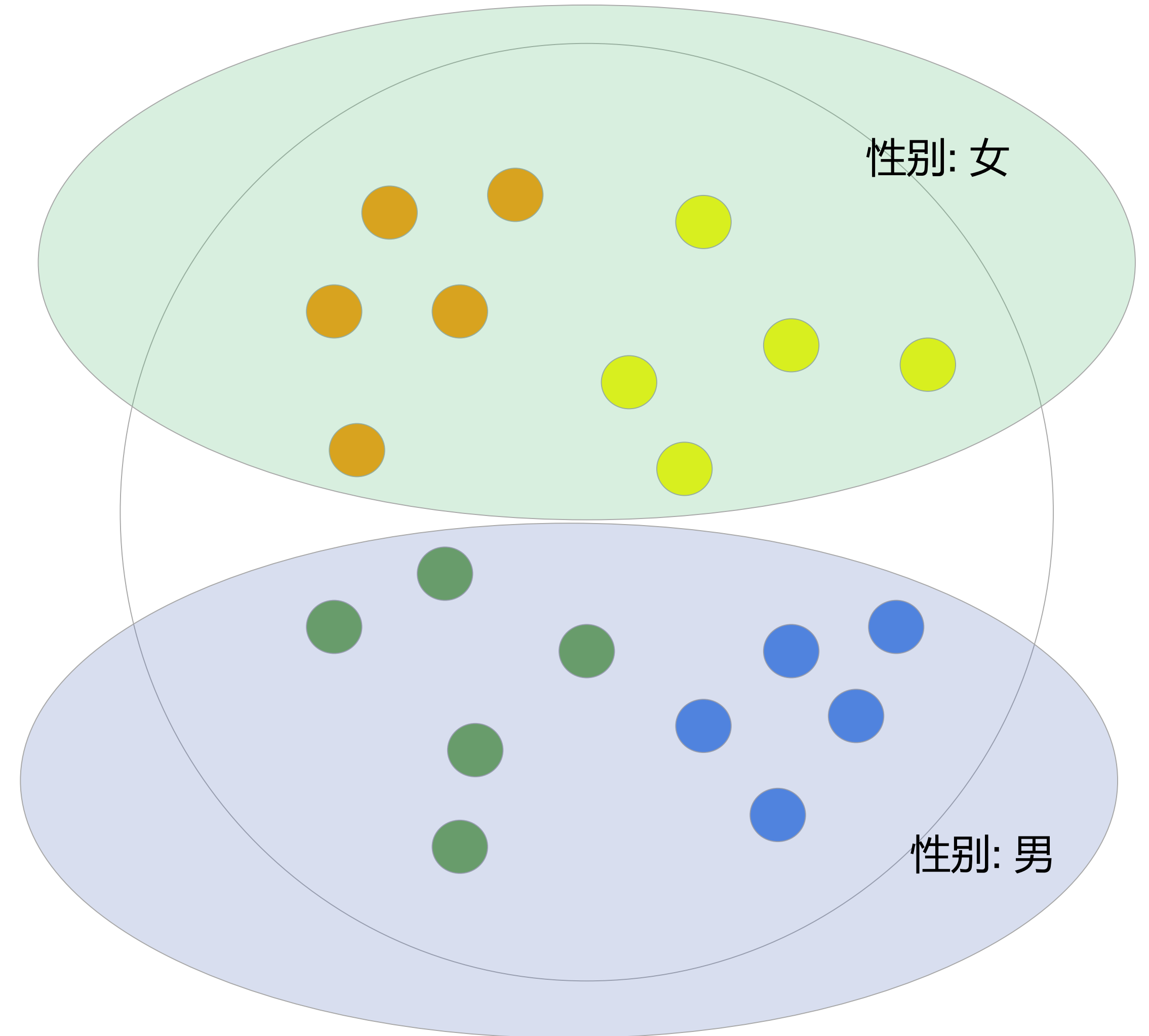
Embedding

Embedding构建-混合标签

通过Word2Vec/GloVe 可以让相似的广告 聚类在一起，包含一定的语义信息

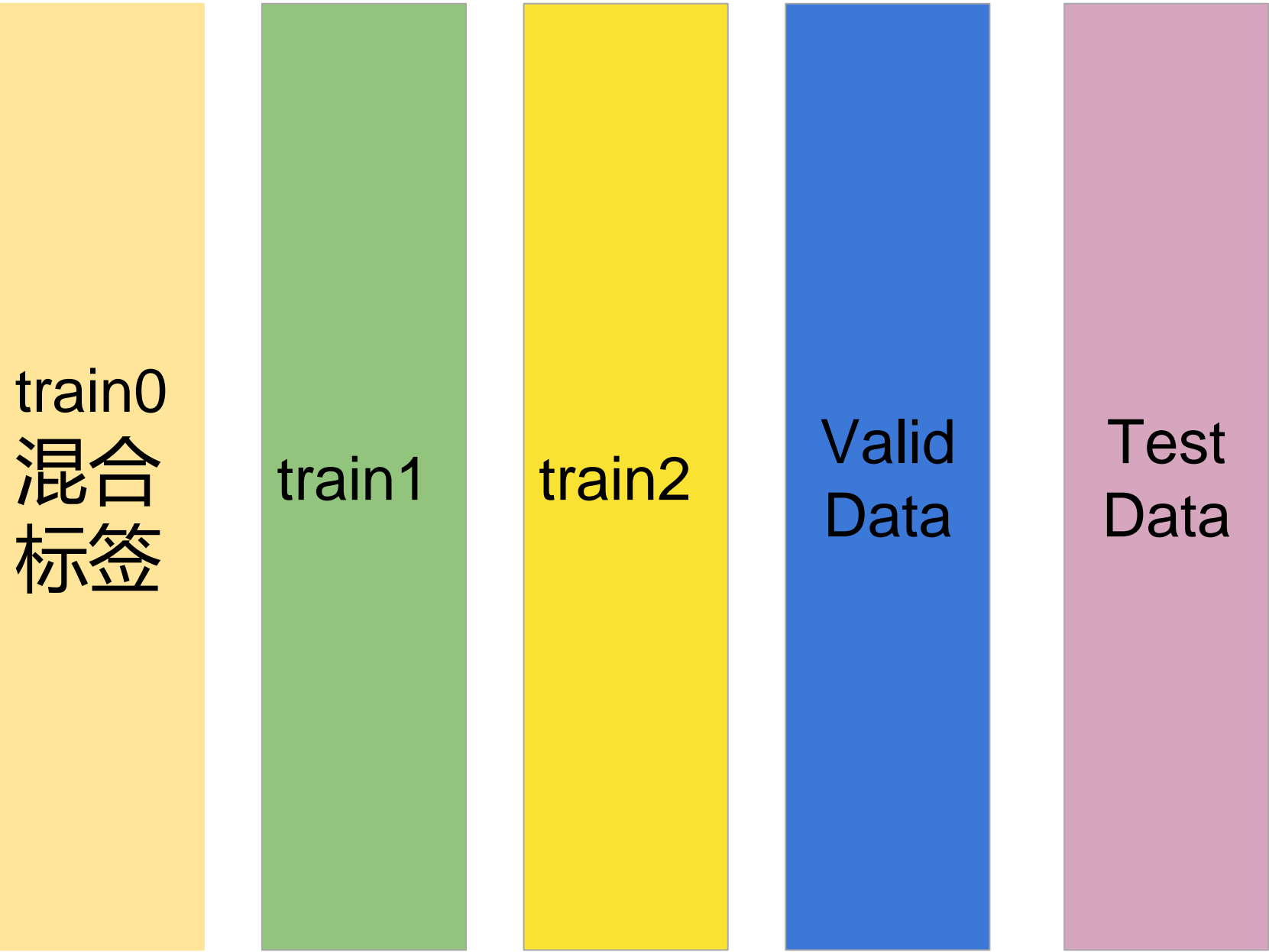
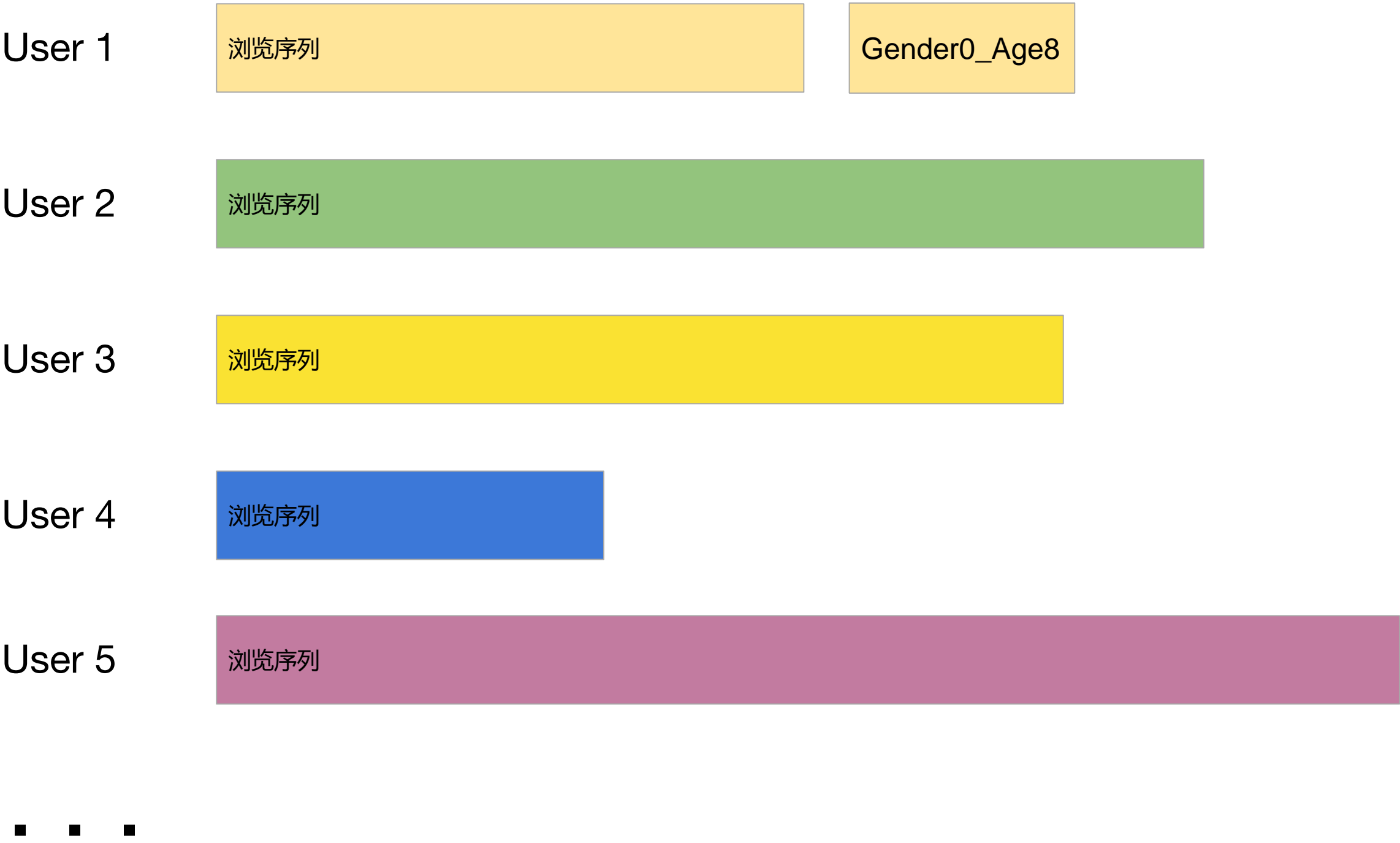


我们希望预训练的Embedding表达语义对下游分类任务更友好



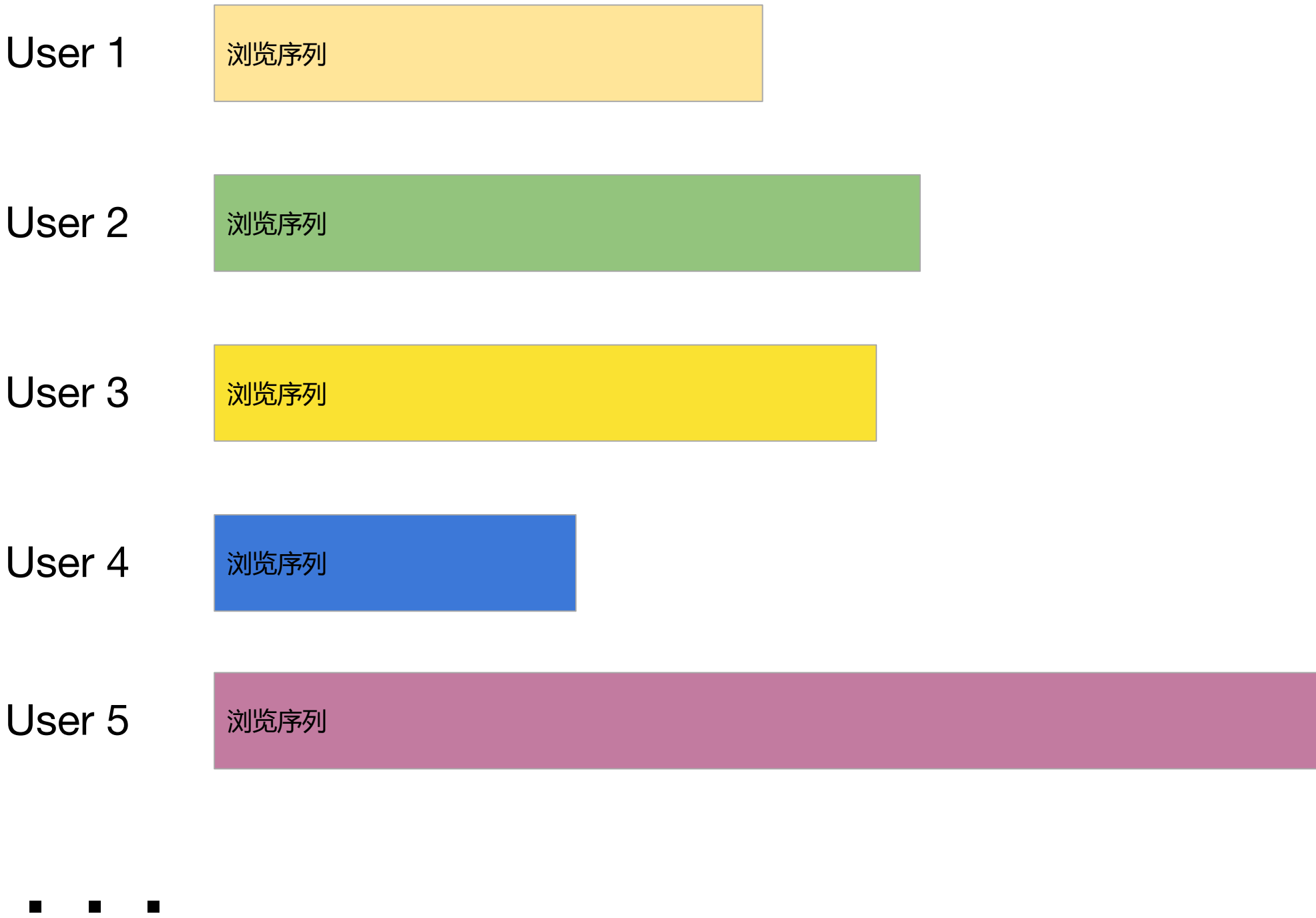
Embedding构建-混合标签

1/4的训练数据混合标签信息训练Word2Vec Embedding

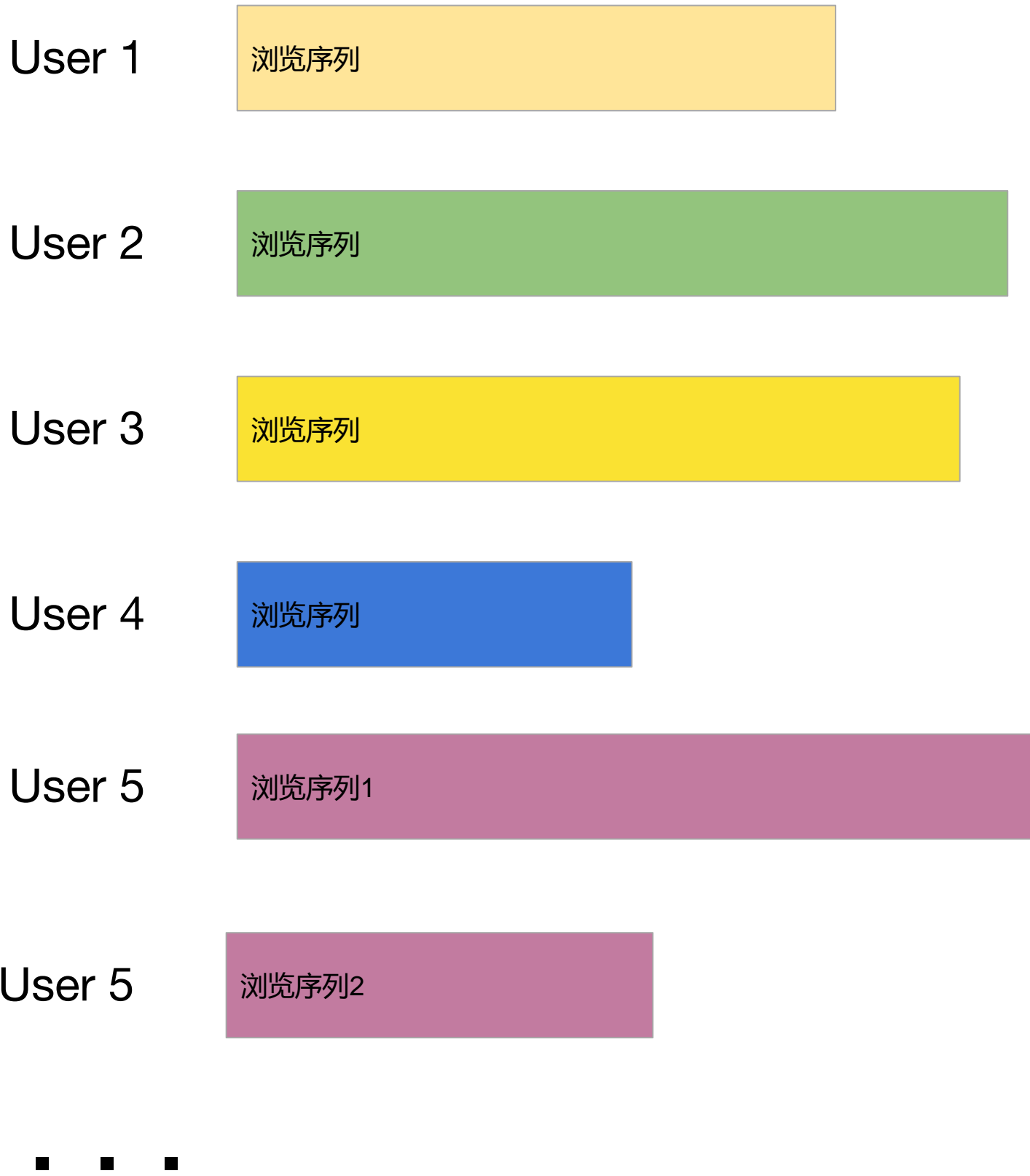


拆分超长浏览序列

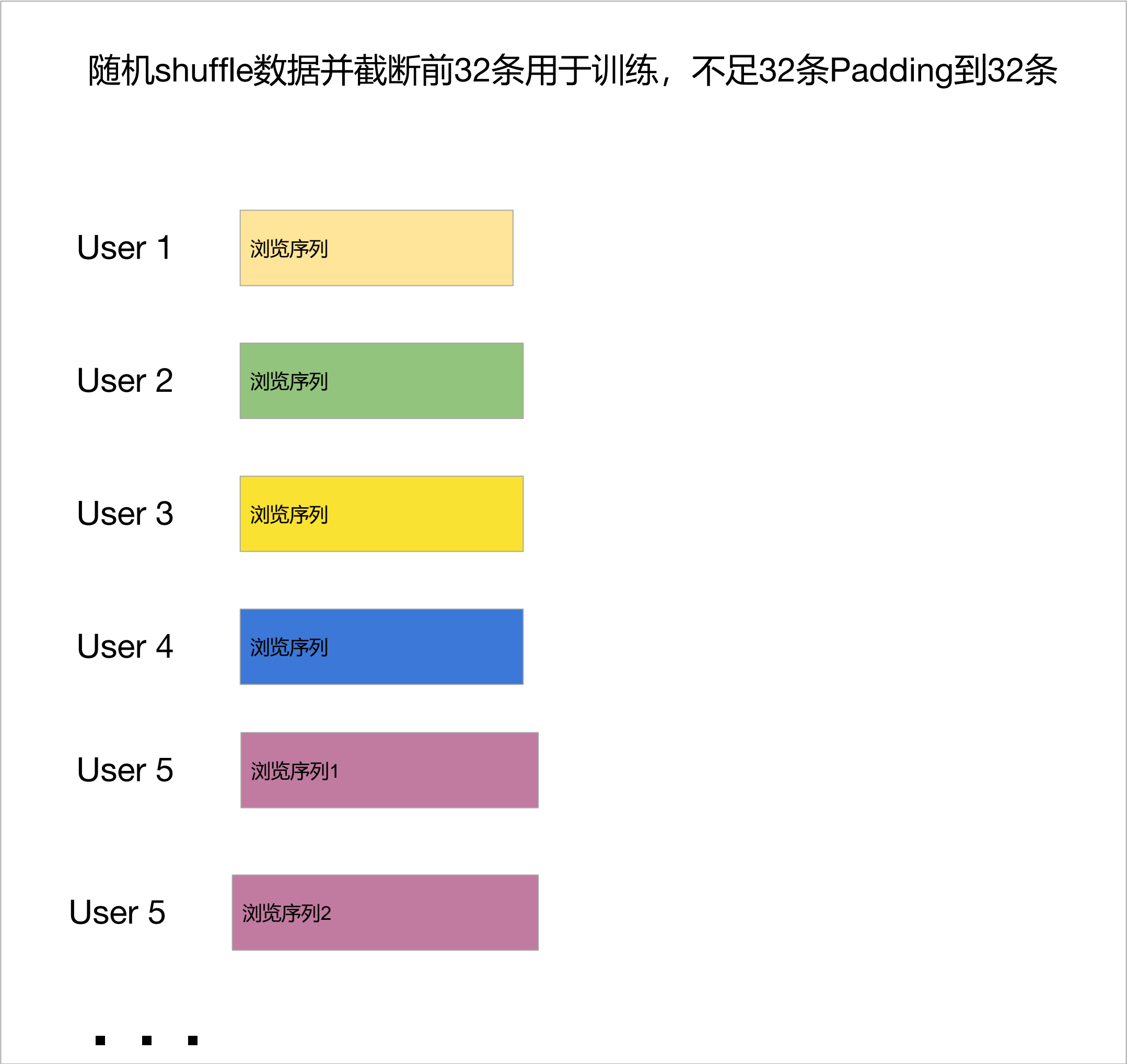
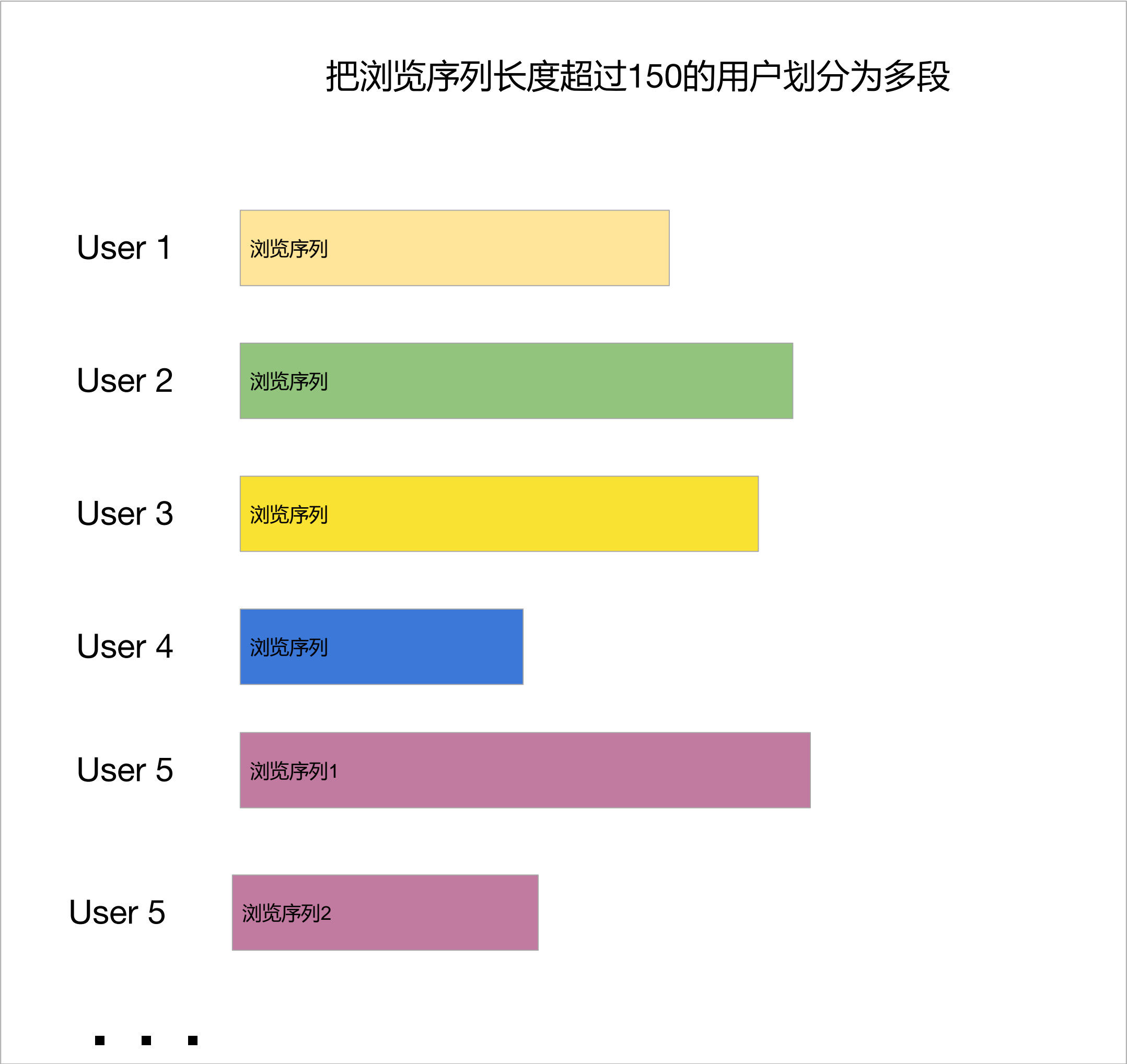
通过点击日志构造每个用户的浏览序列



把浏览序列长度超过150的用户划分为多段



在线数据增广：shuffle 并截断数据

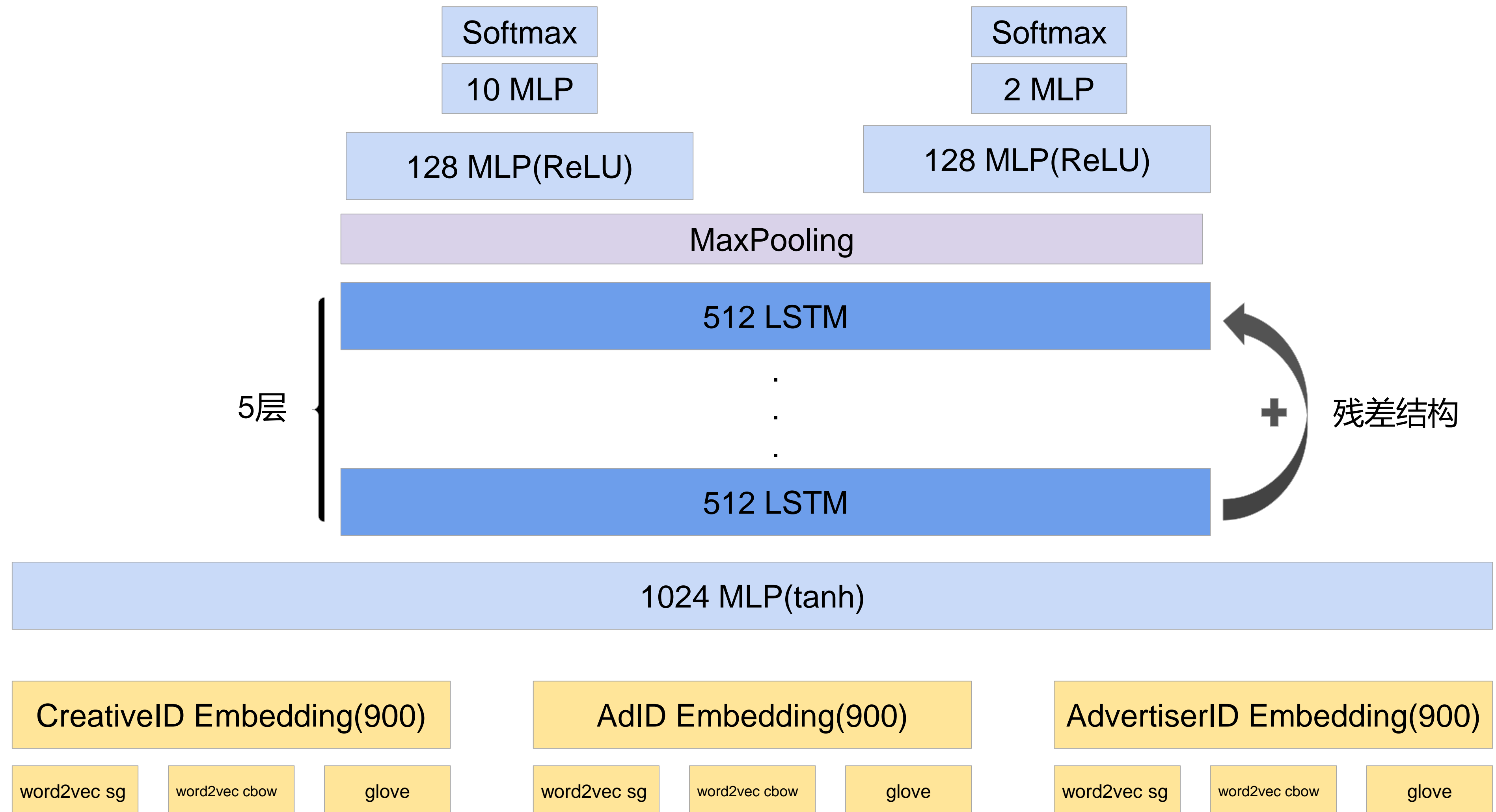




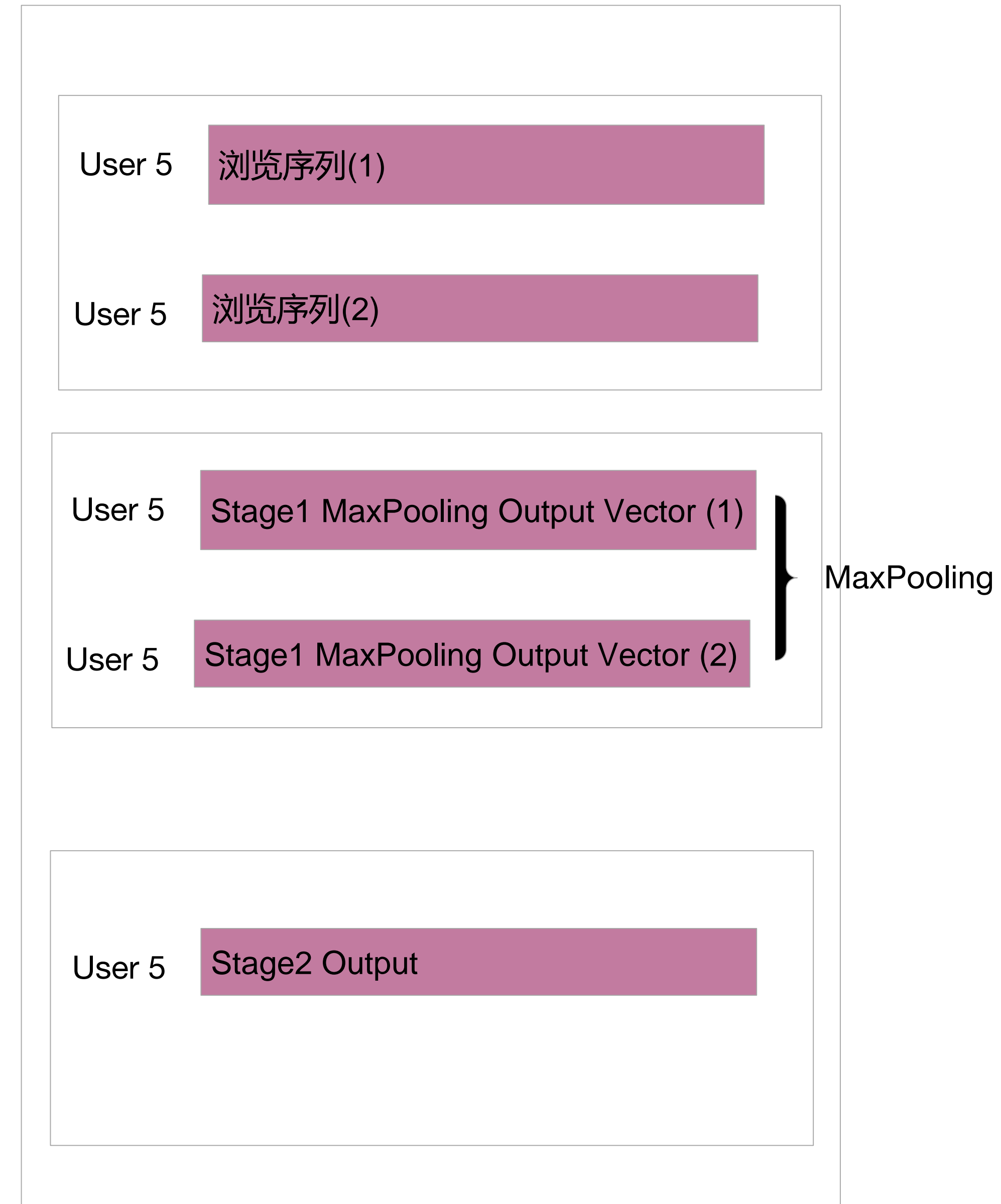
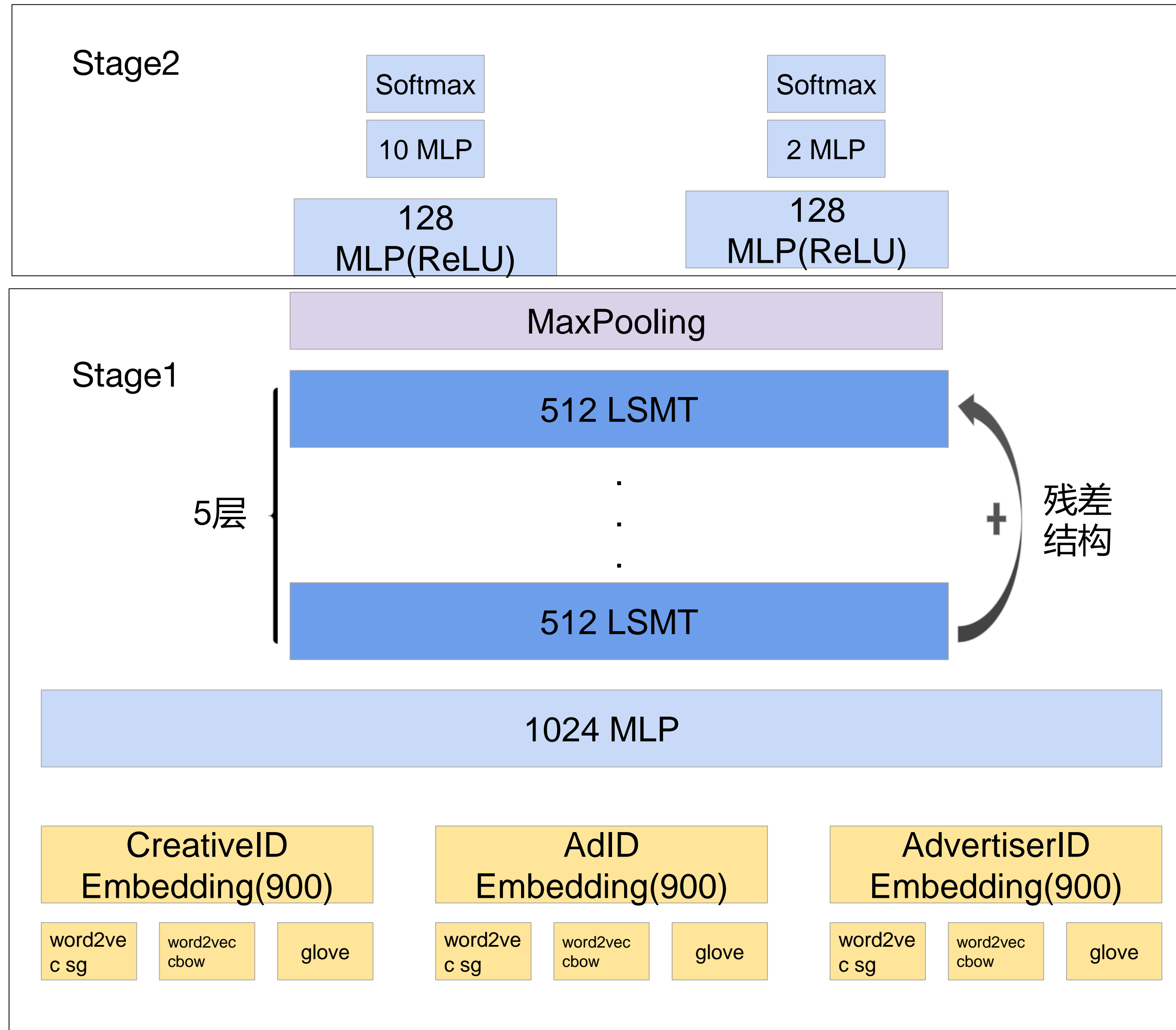
模型介绍

模型结构/模型训练/模型预测

模型结构 - 既要宽更要深



模型结构 - 分段处理超长序列



模型训练 - 天下武功唯快不破

总共400万训练数据, 200S 一个Epoch, 模型收敛到(线下age 0.52, gender 0.95) 需要30个epoch, 共1.7个小时

1. 高效的数据输入, TFRecord +
tf.lookup.StaticHashTable Join 广告特征

2. 训练数据随机shuffle并截取前32条

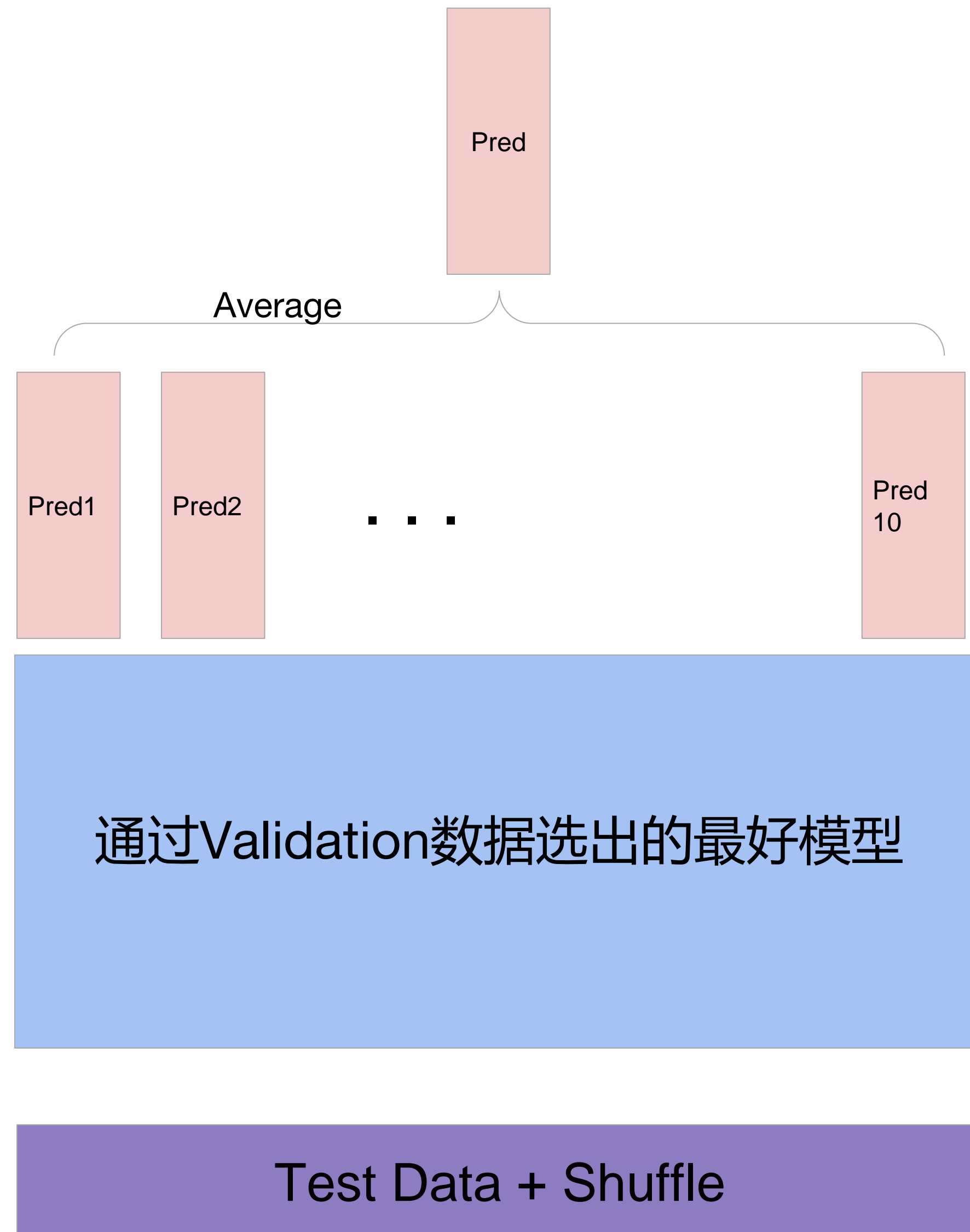
对比150序列长度, 无精度损失并提速4X

3. 过滤出现次数少于3次的ID, 使用FP16存储
Embedding, 使得模型可以放进显存

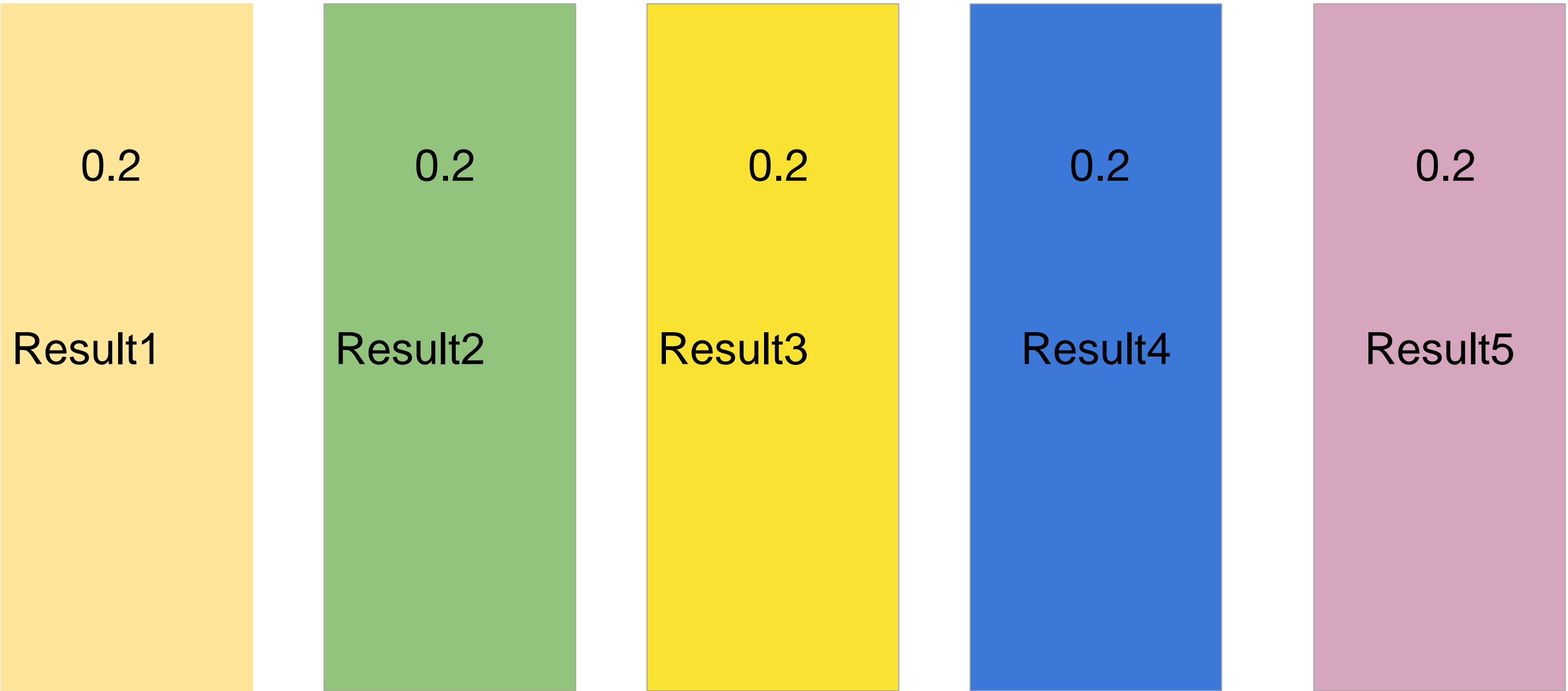
对比Embedding放在CPU, 提速1.3X

4. 开启混合精度训练, 使用TensorCore的能力

对比FP32训练, 提速2.5X



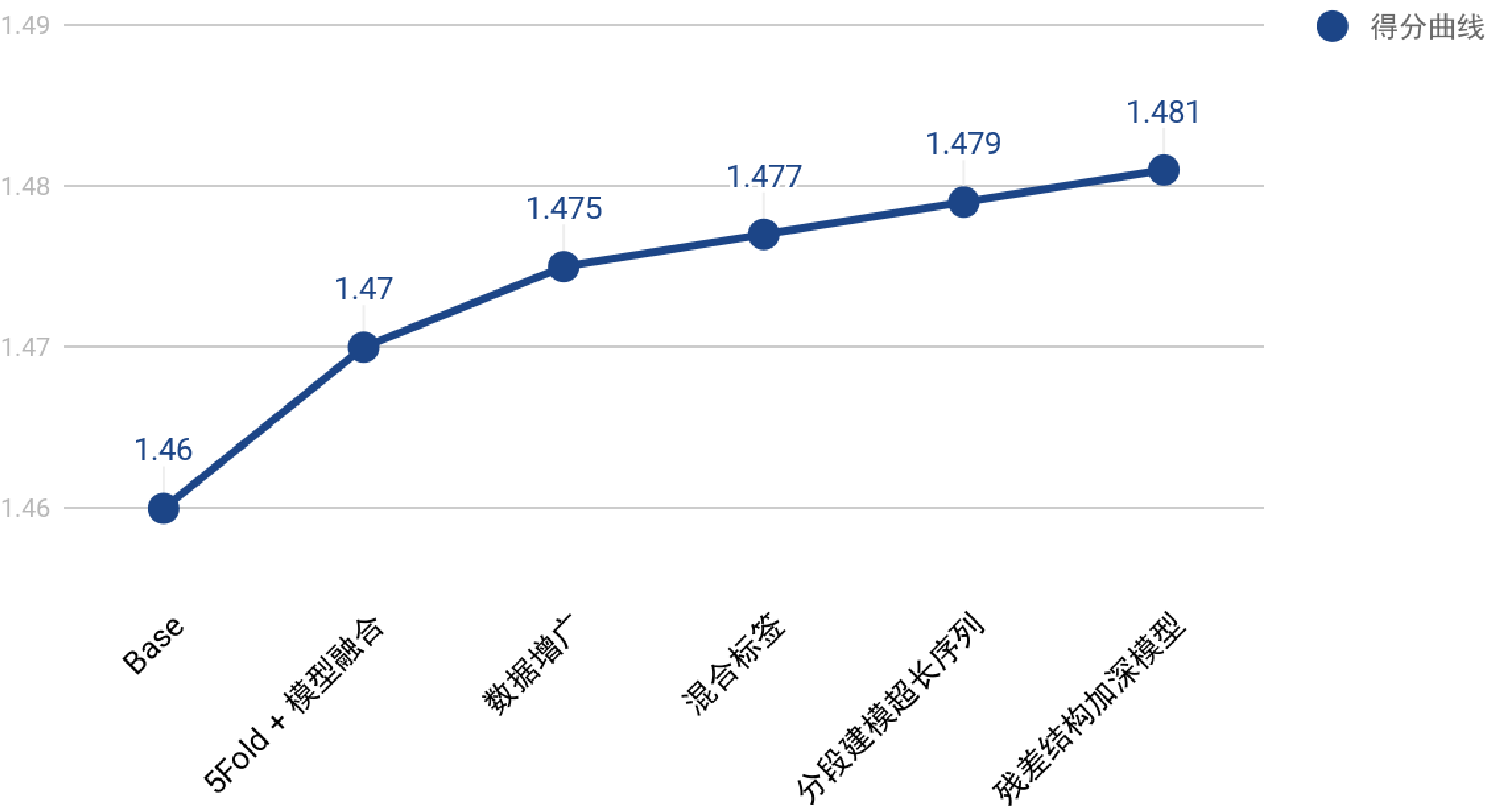
不同的模型最后加权平均得到最终结果





总结与思考

模型得分



总结与思考

主要创新

1. 提出了混合标签构建Word2Vec训练数据，训练更适合下游任务的Embedding
1. 提出了分阶段建模解决数据中的超长序列，训练数据随机采样到32条，在显著提高训练速度的同时不丢失模型精度

问题思考

1. 如何通过CTR日志得到优秀的Embedding表达是解题的关键，除了Word2Vec/GloVe等方式外使用更为复杂的FM/Wide&Deep等CTR模型构建AD和用户的Embedding表达或许可以进一步的提高结果。
1. 我们一直没有利用到时序特征，如何利用上时序特征是另一个可能的得分点

THANK
S

