

陈英发

chenyingfa1999@qq.com | 188 0117 9013 | www.github.com/chen-yingfa

北京 | 研究方向: NLP、LLM、长文本、知识更新



教育经历 EDUCATION

清华大学 博士 | 计算机科学与技术, 自然语言处理

- GPA: 3.9/4.0
- THUNLP 实验室, 研究方向: 长文本建模、知识更新

2022 年 9 月 - 现在
北京

清华大学 本科 | 计算机科学与技术

- GPA: 3.4/4.0
- 二学位: 数字媒体设计

2018 年 8 月 - 2022 年 7 月
北京

发表文章 PUBLICATIONS

Multi-Modal Multi-Granularity Tokenizer for Chu Bamboo Slip Scripts (under review)

Yingfa Chen et al.

- 构造首个楚简（春秋时代的文字）的数据集，含有超过 100K 字和对应的标签。
- 设计一种将字拆解为部件的方式，方便处理没有对应的现代汉字的楚简字。

Beyond the Turn-Based Game: Duplex Models Enable Real-Time Conversations (under review)

Xinrong Zhang, Yingfa Chen et al.

- 提出“双工模型”概念，打破回合制对话格式，模型需要自己判断什么时候开始和停止说话。
- 通过将回合制格式的对话切分，让现有大模型模拟双工的对话模式。

Robust and Scalable Model Editing for Large Language Models (COLING 2024)

Yingfa Chen et al.

- 提出一种基于检索的大模型知识更新的框架，可以同时处理串行和并行的更新操作。
- 构造一个更具有挑战性的知识更新评测数据集。

∞ -Bench: Extending Long Context Evaluation Beyond 100K Tokens (ACL 2024)

Xinrong Zhang, Yingfa Chen et al.

- 构造首个超过 100K 长度的大模型评测集。
- 包含不同语言（中英）和不同领域（数学，代码，自然文本），同时包含合成任务和真实任务。

CFDBench: A Large-Scale Benchmark for Machine Learning Methods in Fluid Dynamics (preprint)

Yining Luo, Yingfa Chen et al.

- 构造了首个针对机器学习模型的，包含多种边界条件、几何形状和流体物性的流体力学评测集。

Sub-Character Tokenization for Chinese Pretrained Language Models (TACL 2023)

Yingfa Chen, Chenglei Si, Zhengyan Zhang et al.

- 将汉字根据字形或者发音转换为 sub-character 序列。
- 此分词器可以获得更短的编码从而提高训练速度，以及对同音字导致的噪声更鲁棒，且它没有牺牲准确率。

READIN: A Chinese Multi-Task Benchmark with Realistic and Diverse Input Noises (ACL 2023)

Yingfa Chen, Chenglei Si, Zhengyan Zhang et al.

BMCook: A Task-agnostic Compression Toolkit for Big Models (EMNLP 2022 Demo)

Zhengyan Zhang, Baitao Gong, Yingfa Chen et al.

其他

- 编程: PyTorch、Python、Huggingface、BMTrain、C++。
- 语言: 普通话、粤语、英语（流利）、挪威语（近乎母语）。
- 兴趣: 羽毛球（系队）。
- 竞赛经历: 数学和信息学（挪威国家队）。
- 个人背景: 本人是挪威籍三代华裔，父母分别是越南和柬埔寨华裔，本人在挪威出生。