# 大语言模型的异构计算和加速

戴金权 (Jason Dai)

英特尔院士

人工智能产业链联盟

星主： AI产业链盟主

○ 知识星球

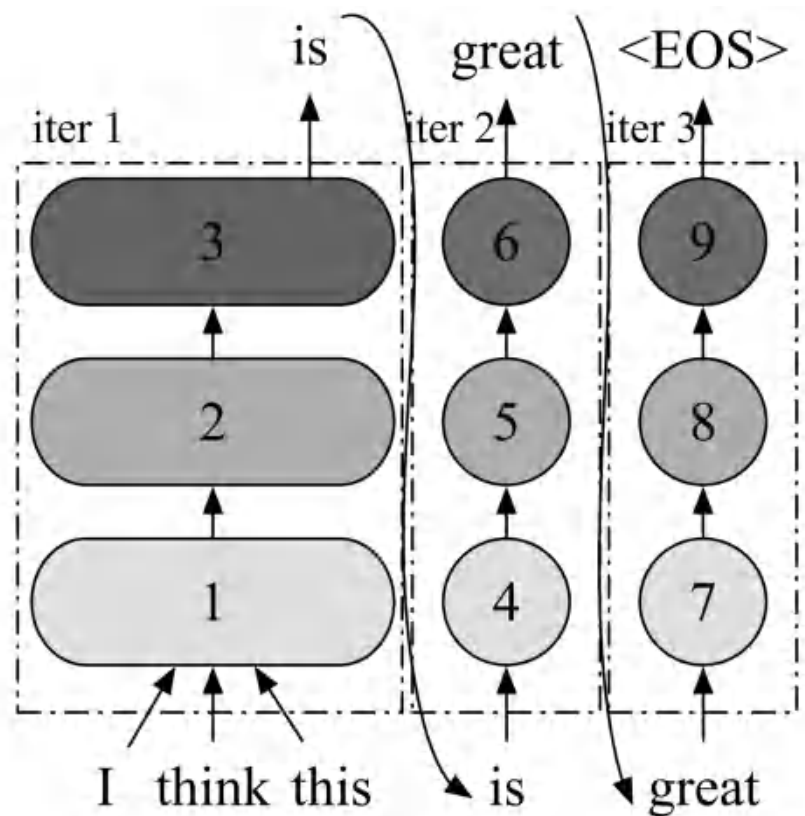微信扫描预览星球详情

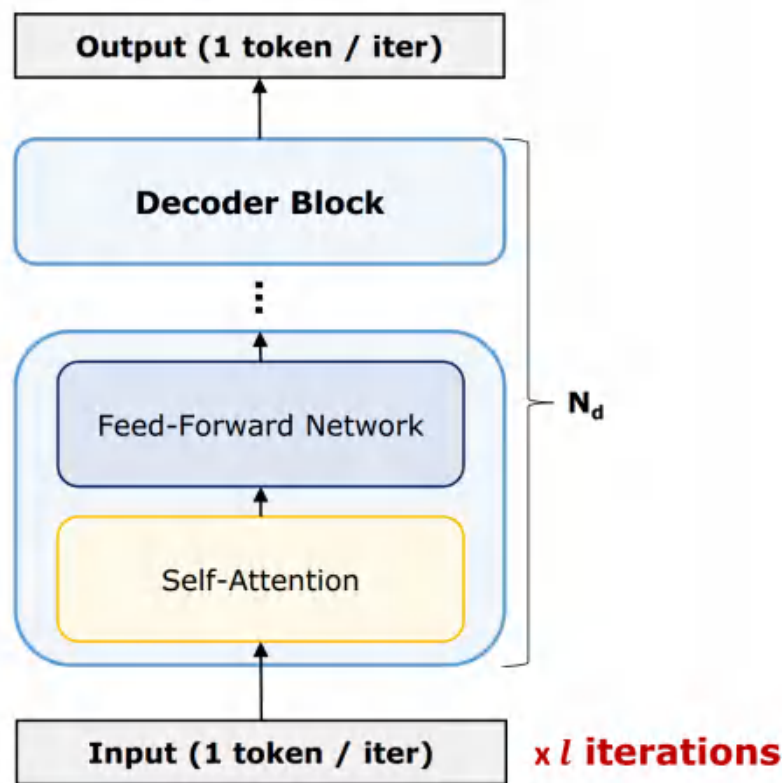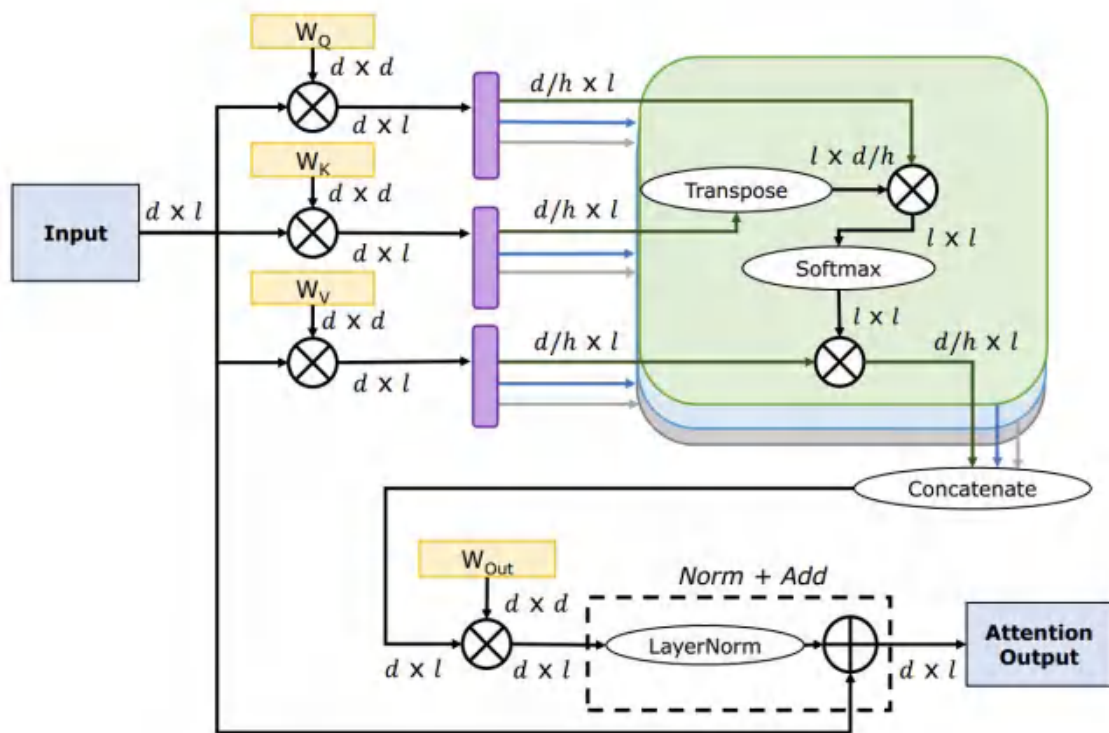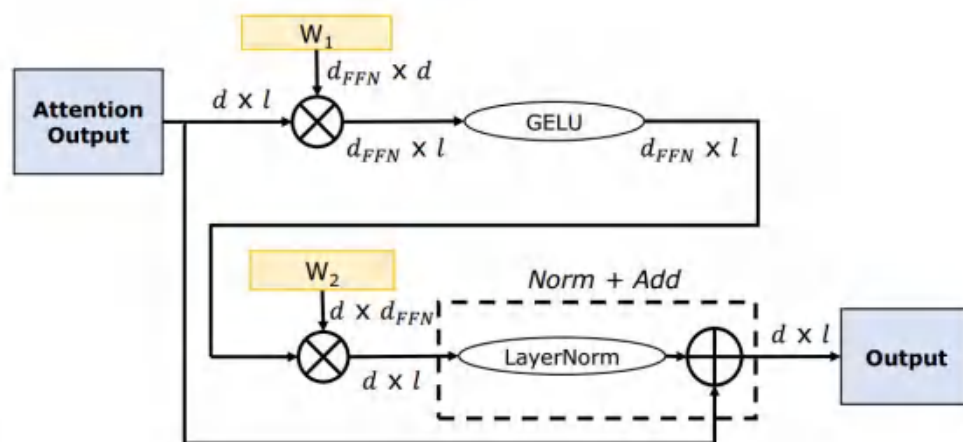# 自回归大语言模型(基于Transformer解码器架构 )



自回归大语言模型：预测下一个token



Transformer解码器架构

# Transformer解码器架构



**Muti-Head Attention (MHA) Module**

**Feed-Forward Network (FFN) Module**

训练; 推理 (第一个token/Prefill)
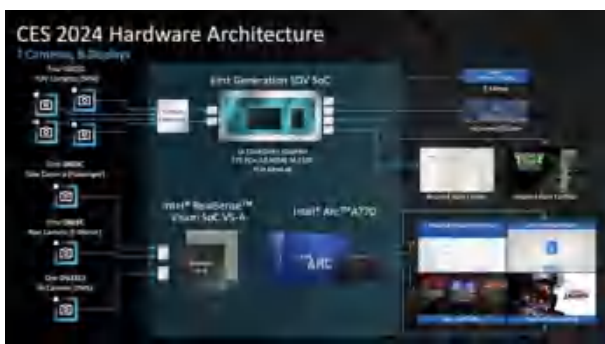
# Transformer解码器架构



推理 (下一个token/Decode)

# 大语言模型推理和训练瓶颈

- 内存带宽

- 计算

- 显存大小

- 分布式计算 (互联)

# 大模型的异构计算和加速

- XPU异构计算
  - CPU, GPU, NPU硬件加速



客户端
（ Intel Core Ultra AI PC ）



边缘端
（ Intel AI座舱 ）



2~10 x Arc A770 GPU (16GB)

服务器
（ Intel Xeon+Intel Arc GPUs ）

# 大模型的异构计算和加速

▪ 低比特计算

- 模型量化/压缩 (WxAy)

- 数据类型 (INTx, FPx)

- 低比特算子

- 显存(如kv cache) 使用量

- 训练、微调 (如QLoRA)

# 低比特大模型的精度

困惑度 (Wikitext数据集)

| Perplexity | sym_int4 | q4_k | fp6 | fp8_e5m2 | fp8_e4m3 | fp16 |
|---|---|---|---|---|---|---|
| Llama-2-7B-chat-hf | 6.364 | 6.218 | 6.092 | 6.180 | 6.098 | 6.096 |
| Mistral-7B-Instruct-v0.2 | 5.365 | 5.320 | 5.270 | 5.273 | 5.246 | 5.244 |
| Baichuan2-7B-chat | 6.734 | 6.727 | 6.527 | 6.539 | 6.488 | 6.508 |
| Qwen1.5-7B-chat | 8.865 | 8.816 | 8.557 | 8.846 | 8.530 | 8.607 |
| Llama-3.1-8B-Instruct | 6.705 | 6.566 | 6.338 | 6.383 | 6.325 | 6.267 |
| gemma-2-9b-it | 7.541 | 7.412 | 7.269 | 7.380 | 7.268 | 7.270 |
| Baichuan2-13B-Chat | 6.313 | 6.160 | 6.070 | 6.145 | 6.086 | 6.031 |
| Llama-2-13b-chat-hf | 5.449 | 5.422 | 5.341 | 5.384 | 5.332 | 5.329 |
| Qwen1.5-14B-Chat | 7.529 | 7.520 | 7.367 | 7.504 | 7.297 | 7.334 |

# 大模型的异构计算和加速

▪ 推理算法优化

- Self-speculative decoding

- KV Cache compression

- Sliding window attention

- Sparse attention

- Flash attention/decoding

- Continuous batching

- Prefill/decoding disaggregation

- …

# IPEX-LLM: 开源大模型XPU加速框架

Users/Developers

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

**Python (PyTorch) Ecosystem**

HuggingFace,
Langchain,
LlamaIndex,
DeepSpeed,
TRL, Axolotl,
...

**llama.cpp Ecosystem**

llama.cpp,
Ollama,
LangChain.js,
Open WebUI,
...

**IPEX-LLM Library**

**XPU Compute**

**LLM Acceleration**

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

Intel XPU

*https://github.com/intel-analytics/ipex-llm/*

# 英特尔 XPU 大模型加速体验

# Intel UHD/Iris iGPU
llama.cpp + IPEX-LLM (Phi-3-mini, Q4_0)

# Intel Core Ultra AI PC
Ollama + IPEX-LLM (Mistral-7B, Q4_K_M)

# Intel Arc A770 GPU

TextGeneration-WebUI + IPEX-LLM (Llama3-8B, FP8)

# 4 x Arc A770 GPU
FastChat + IPEX-LLM (QWen1.5-72B FP6)

# LoRA/QLoRA on Xeon+Multi-Arc
支持 PEFT, TRL, Axolotl, Zero2/Zero3

# 英特尔 XPU 大模型应用创新

# Office助手

ExtendOffice展示

# 工业机器人代码生成

科东软件展示



转换的结果，直接控制机器人手臂进行操作

# AI座舱-汽车助理

智谱AI展示

# AI座舱-驾驶伴侣

百川智能展示

# 个人或企业本地RAG系统



在英特尔 XPU 上运行 RAGFlow
*(https://github.com/intel-analytics/ipex-llm/blob/main/docs/mddocs/Quickstart/ragflow_quickstart.md)*

在英特尔 XPU 上运行 GraphRAG
*(https://github.com/intel-analytics/ipex-llm/blob/main/docs/mddocs/Quickstart/graphrag_quickstart.md)*

# AI人工智能产业链联盟

## #每日为你摘取最重要的商业新闻#

### 更新 · 更快 · 更精彩

> **Zero**
> AI音乐创作人
> 水墨动漫联盟创始人
> 百脑共创联合创始人
> 人工智能产业链联盟创始人
> 中关村人才协会秘书长助理
> 河北北大企业家分会秘书长
> 器玫星辰智能科技有限公司CEO
> 河北清华发展研究院智能机器人中心线上负责人
> 中关村人才协会数字体育与电子竞技专委会秘书长助理

> **主要业务:**AI商业化答疑及课程应用场景探索，各类AI产品学习手册，答疑及课程

**欢迎扫码交流**

提供：学习手册/工具/资源链接/商业化案例/
行业报告/行业最新资讯及动态

---

**人工智能产业链联盟创始人**
邀请你加入星球，一起学习

## 人工智能产业链联盟报告库

星主：人工智能产业链联盟创始人

每天仅需0.5元，即可拥有以下福利！
每周更新各类机构的最新研究成果。立志将人工智能产业链联盟打造成市面上最全的AI研究资料库，覆盖券商、产业公司、研究院所等...

**知识星球**

微信扫码加入星球 ▶

# Call to Actions

- 关注和试用 IPEX-LLM，并给我们反馈
  - *https://github.com/intel-analytics/ipex-llm/*


- 使用 IPEX-LLM 在 Intel XPU 平台开发大模型及其应用
  - 客户端-边缘-服务器（Intel Core Ultra AI PC、AI座舱、Xeon+Intel Arc GPUs）
  - 高效的大模型 XPU 加速的创新
  - 大模型应用场景的创新

谢谢!

# Notices &Disclaimers

Performance varies by use, configuration and other factors. Learn more on the Performance Index site.

Performance results are based on testing as of dates shown in configurations and may not reflect all publicly available updates.  See backup for configuration details.  No product or component can be absolutely secure.

Your costs and results may vary.

Intel technologies may require enabled hardware, software or service activation.

© Intel Corporation.  Intel, the Intel logo, and other Intel marks are trademarks of Intel Corporation or its subsidiaries.  Other names and brands may be claimed as the property of others.