# 华为盘古大模型的核心技术与挑战

## 刘群 LIU Qun

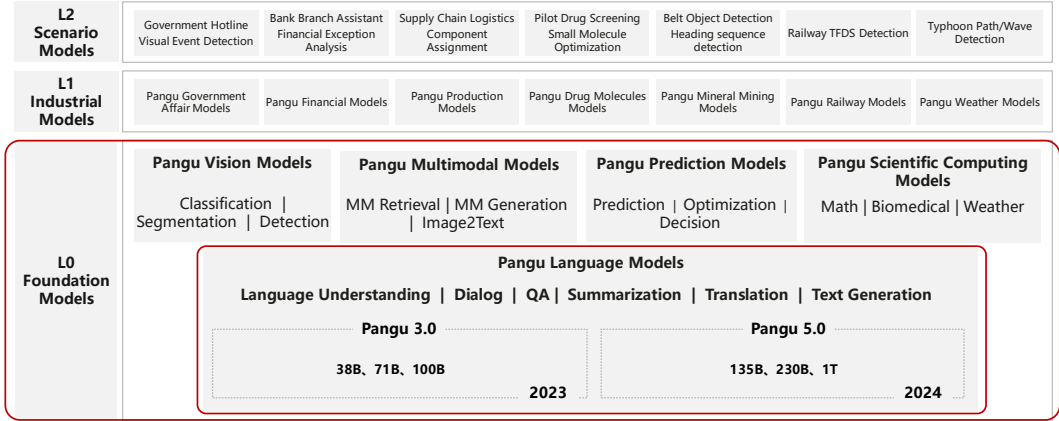**Huawei Noah's Ark Lab**

**CCF大模型论坛**

**2024.06.06, 北京**

NOAH'S ARK LAB

HUAWEI

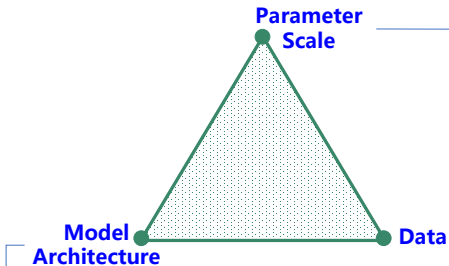# Pangu Large Models: AI4Industry

► Huawei regards AI as a huge and crutial oppotunity for the future of the company.
► Huawei has invested the full-stack AI technologies for AI, including NPU chips (Ascend), clusters (Atlas), AI frameworks (MindSpore), AI models (Pangu), and a broad spectrum of AI applications, especially for industries.

| L2 Scenario Models | Government Hotline Visual Event Detection | Bank Branch Assistant Financial Exception Analysis | Supply Chain Logistics Component Assignment | Pilot Drug Screening Small Molecule Optimization | Belt Object Detection Heading sequence detection | Railway TFDS Detection | Typhoon Path/Wave Detection |
|---|---|---|---|---|---|---|---|

| L1 Industrial Models | Pangu Government Affair Models | Pangu Financial Models | Pangu Production Models | Pangu Drug Molecules Models | Pangu Mineral Mining Models | Pangu Railway Models | Pangu Weather Models |
|---|---|---|---|---|---|---|---|

**L0 Foundation Models**

| **Pangu Vision Models** | **Pangu Multimodal Models** | **Pangu Prediction Models** | **Pangu Scientific Computing Models** |
|---|---|---|---|
| Classification \| Segmentation \| Detection | MM Retrieval \| MM Generation \| Image2Text | Prediction \| Optimization \| Decision | Math \| Biomedical \| Weather |

**Pangu Language Models**

**Language Understanding | Dialog | QA | Summarization | Translation | Text Generation**

| **Pangu 3.0** | **Pangu 5.0** |
|---|---|
| **38B、71B、100B** | **135B、230B、1T** |
| 2023 | 2024 |

# Beyond Algorithms: Navigating the Data Deluge in AI

▶ Building large-scale AI models has become a massive systems engineering problem, far more than just an algorithm problem, which requires cooperations among scientists and engineers from multiple disciplines.

**Parameter Scale**

**Model Architecture**

**Data**

The more parameters, the more intelligence, requiring more computing power.
- **Training:**
  - **Large-scale cluster computing power:** Huawei Atlas 900 clusters are capable of training models at the scale of trillion parameters.

    *Atlas900 SuperCluster*

  - **Memory reduction:** Training larger models in given clusters.
- **Inference:** ultimate quantization compression with almost lossless precision

Transformer is the current mainstream model architecture.
Extensions:
- **Sparse-dense architecture (MoE):** high scalability and supporting larger models with less computing power
- **Vector database (RAG):** "hippocampus" for LLMs, external memory
- **Plug-ins and tools:** invoking external tools to complete complex tasks

Data is the source of knowledge and intelligence.
Both "quality" and "quantity" of data are crucial.
- **Pre-training data:** diversity, coverage, cleanness, consistency, and timeliness
- **Instruction fine-tuning data:** instruction following and alignment with human intentions

HUAWEI  NOAH'S ARK LAB

# Data Management: acquisition, cleaning, labeling, and pre-processing

**Pretraining Data ( 10T+ Tokens )**

Papers Codes
Books
Webpages

| Asset Dashboard | Data Labels | Copyright management | Permission control | Security compliance |

## Labeling

| Manual labeling | Reverse labeling | Self-supervised labeling |

## Cleaning

| Data cleaning | Model cleaning | Operator cleaning |

**Acquisition (1,000+TB)**

| **Open web data** web pages and open-source codes | **Purchased data** books, question banks, patents, and periodicals | **Ecosystem cooperation** Industrial open data |

HUAWEI  NOAH'S ARK LAB

# PanGu-Σ : dense-sparse architecture with heterogeneous computing



PanGu-Σ : Towards Trillion Parameter Language Model with Sparse Heterogeneous Computing, arXiv:2303.10845

# CAME: confidence-guided adaptive memory efficient optimization

## 1~4x memory usage reduction, 100% throughput improving

**Challenges in inference: Large number of parameters, slow inference,
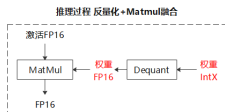high memory usage, high cost in end-to-end inference**

➤ **Traditional quantization causes significant precision degradation for generative models**
➤ **High memory usage in inference: (1) Model parameters: 350 GB memory for a 175B model.
2) KV cache: 576 GB for a 175B model with a 4 KB length.**



### Low-bit Weighting Algorithm: QuantGPT

**Progress:** (1) **4/8-bit quantization** algorithm (2) Ascend affinity efficient **dequantization** operator (3) **2~4x** memory reduction, **15-30%** inference acceleration.
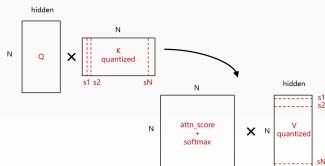**Deployment of a 38B model on a single card.**

推理过程 反量化+Matmul融合

激活FP16 → MatMul ← Dequant ← 权重 IntX
权重 FP16
↓
FP16

arXiv:2203.10705v2 **ACL 2022 Outstanding Paper Awards**

### KV cache Compression

**Progress: 1x** memory reduction after KV cache 8-bit quantization



### Separate deployment, Dynamic batch

**Progress: Separate deployment of** full and incremental inference (8+8), **30-50%** throughput improvement

**Dynamic batch:**

• Early exit for completed samples
• New samples added in time
**Separate deploy for full and incremental inference :**
• Full inference: batch **size**=1 to reduce delay
• Incremental inference: Large Batch Size to improve throughput

Context (用户query+对话历史)
↓
全量推理 (context计算)  输入n个token 输出1个token
↓
Context部分kv cache
↓
增量推理 (逐token生成)  输入1个token 输出1个token
↓
流式解码返回答案

HUAWEI    NOAH'S ARK LAB

# Efficient Post-Training Pruning Method for LLMs

## ① Background

### Typical types of model pruning

- ☐ Unstructured Pruning ✓
- ☐ Structured Pruning ✗
- ☐ Semi-structured Pruning ✓ (2:4 Sparsity)

### Advantages of LLM pruning

| | Memory Access | | Computation | |
|---|---|---|---|---|
| | Model Size | Throughput | Prefill | Decode |
| Unstructured Pruning | ↓ | ↑ | | |
| Structured Pruning | ↓ | ↑ | ↓ | ↓ |
| Semi-structured Pruning | ↓ | ↑ | ↓ | ↓ |

### Challenges in LLM pruning: Channel Collapse

Certain channels exhibit smaller magnitudes, that **500+ channels are pruned with Wanda**

**Nearly no channel collapse with RIA**

input channel degree distribution after pruning

## ② Method

### A new pruning metric



Pruned (50% sparsity) $RI_{2,4} = \frac{|W_{24}|}{\sum|W_{*4}|} + \frac{|W_{24}|}{\sum|W_{2*}|} = 0.45$   Preserved (50% sparsity)

- **Relative Importance and Activation (RIA) :**

$$\mathbf{RIA}_{ij} = \mathbf{RI}_{ij} \times (\|\mathbf{X}_i\|_2)^a = \left( \frac{|\mathbf{W}_{ij}|}{\sum|\mathbf{W}_{*j}|} + \frac{|\mathbf{W}_{ij}|}{\sum|\mathbf{W}_{i*}|} \right) \times (\|\mathbf{X}_i\|_2)^a,$$

jointly normalizes the weight in the input and output dimensions, together with activations

- Effectively resolving the channel collapse issue

### Semi-structured pruning (2:4/4:8 sparsity)

- Permute the channels equivalently to find a better 2:4/4:8 sparse pattern
- **Fast permutation: 15.3s** for a single linear layer, 100x faster than previous permutation method

## ③ Result

### Unstructured Pruning

Table 1: Perplexity results on Wikitext2. We produce the one-shot Post-Training pruning methods with 50% unstructured sparsity on LLaMA, LLaMA2 and OPT models.

| Method | LLaMA 7b | LLaMA 13b | LLaMA 30b | LLaMA 65b | LLaMA2 7b | LLaMA2 13b | LLaMA2 70b | OPT 1.3b | OPT 13b |
|---|---|---|---|---|---|---|---|---|---|
| Dense | 5.68 | 5.09 | 4.77 | 3.56 | 5.47 | 4.88 | 3.32 | 14.62 | 10.13 |
| Magnitude | 17.28 | 20.22 | 7.54 | 5.90 | 16.02 | 6.83 | 5.36 | 1712 | 11561 |
| Wanda | 7.26 | 6.15 | 5.24 | 4.57 | 6.92 | 5.99 | 4.22 | 18.41 | 11.92 |
| SparseGPT | 7.24 | 6.20 | 5.32 | 4.57 | 6.99 | 6.10 | 4.25 | 27.00 | 11.18 |
| RIA (Ours) | **7.12** | **6.08** | **5.08** | **4.38** | **6.81** | **5.83** | **4.11** | **18.08** | **11.05** |

- SOTA performance under 50% sparsity for 7B-65B model
- No parameter update and all sizes of LLMs can be compressed within **seconds**

### Semi-structured Pruning (2:4/4:8)

Table 5: LLaMA2-70B: Zero-Shot Performance of N:M constraint model comparing to the dense model. Bold values denote the best performance across all N:M constraint models. An asterisk ("*") signifies performance surpassing that of the dense method.

| Method | Hellaswag | BoolQ | ARC-C | MNLI | RTE | AVG |
|---|---|---|---|---|---|---|
| Dense | 64.77 | 83.70 | 54.44 | 45.81 | 67.87 | 63.32 |
| Wanda (2:4) | 57.35 | 81.44 | 46.01 | 37.69 | 68.59* | 58.22 |
| Wanda (2:4+CP) | 59.37 | 84.50* | 48.55 | 43.09 | 66.43 | 60.39 |
| Wanda (4:8+CP) | 60.86 | 82.73 | 49.94 | 40.15 | 67.87 | 60.31 |
| RIA (2:4) | 57.13 | 82.78 | 46.76 | 37.39 | 69.31* | 58.68 |
| RIA (2:4+CP) | 58.48 | 85.14* | 49.15 | 49.08* | 68.95* | 62.16 |
| RIA (4:8+CP) | 60.44 | 83.58 | 50.43 | 48.69* | 70.04* | 62.64 |

- Only **~1%** average Acc drop
- The channel permutation semi-sparsity only takes **40** minutes (LLaMA 65B)

Zhang, Y., et.al, Plug-and-Play: An Efficient Post-Training Pruning Method for Large Language Models. ICLR 2024    Hall B #225, 4:30PM – 6:30PM

7   total: 14
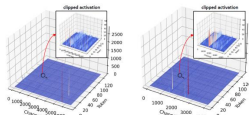
HUAWEI    NOAH'S ARK LAB

# IntactKV: An orthogonal solution to enhance quantized LLMs

## 1 Background

**We discover Pivot tokens:**
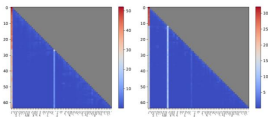(usually the BOS and delimiter tokens ( "/" , "." , "," , ".") )
- Activations with extremely large values



(a) Output activations of LLaMA-30B Layer 24

(b) Output activations of LLaMA-2-7B Layer 24

- Highly concentrated attention scores over these tokens
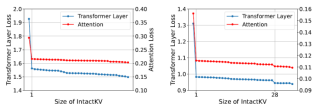- Pivot tokens are sensitive to quantization



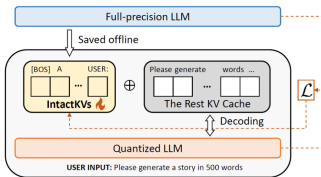(c) Attention map of LLaMA-30B Layer 24

(d) Attention map of LLaMA-2-7B Layer 24

## 2 Method

**Keeping the KV cache of pivot tokens intact (i.e., generating them from the full-precision model) can effectively lower the quantization error**



(a) LLaMA-13B

(b) LLaMA-30B

**System Prompt:** [BOS] A chat between ... intelligence assistant .... USER:



(a) The overview of INTACTKV.

- **IntactKVs** are concatenated with the normal KV cache from user queries

## 3 Result

**Consistent improvement** over existing methods on MMLU benchmark

| Task Acc | MMLU (5 shot) average | | | | |
|---|---|---|---|---|---|
| Vicuna Family | v1.5-7B | v1.5-13B | v1.3-7B | v1.3-13B | v1.3-33B |
| FP16 | 49.84% | 55.78% | 47.12% | 52.10% | 59.30% |
| RTN | 44.62% | 51.44% | 39.33% | 44.56% | 53.18% |
| GPTQ | 43.99% | 52.95% | 40.12% | 47.83% | 55.84% |
| OmniQuant | 46.54% | 52.86% | 43.18% | 47.92% | 55.12% |
| AWQ | 46.45% | 52.92% | 43.08% | 48.56% | 56.09% |
| +INTACTKV[B] | 46.87% | 53.58% | 44.67% | 49.05% | 56.91% |

**Integratable with KV Cache Quantization**



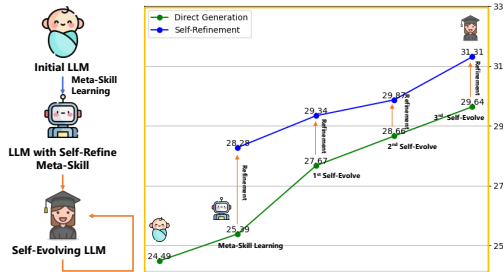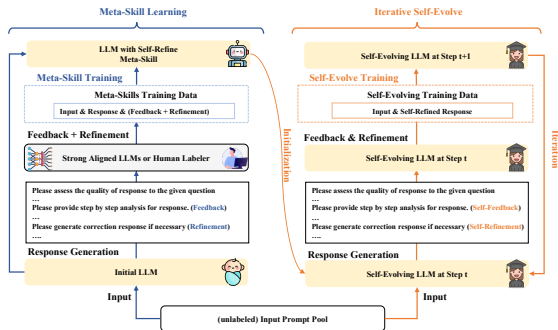(a) Vicuna-v1.3-7B

(b) Vicuna-v1.3-13B

**Other advantages:**
- Orthogonal to existing quantization methods
- Plug-and-play: no extra training /inference overhead

Liu, R., et.al, IntactKV: Improving Large Language Model Quantization by Keeping Pivot Tokens Intact, arXiv: 2403.01241, 2024
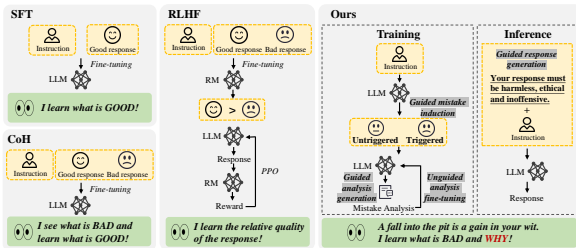
HUAWEI    NOAH'S ARK LAB

# SELF: self-improving and self-evolving for LLMs

- ► Freely available high quality data is going to be exhausted in the near future.
- ► Instruct data for SFT and human preference data for RLHF are expensive.
- ► We introduce an innovative approach, SELF, which empowers LLMs to undergo continual self-evolution, thereby augments their inherent capabilities.



SELF: Self-Evolution with Language Feedback, arXiv:2310.00533v2

# Gaining Wisdom from Setbacks: Aligning LLMs via Mistake Analysis



Gaining Wisdom from Setbacks: Aligning LLMs via Mistake Analysis, ICLR 2024, arXiv:2310.10477v1
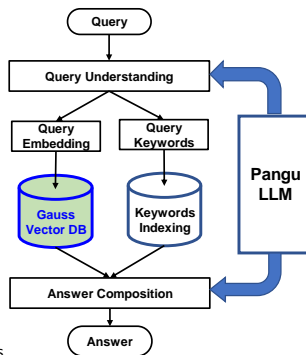
# RAG with Vector Database

## Application Case: Telecom Knowledge Q&A

**Telecom Knowledge QA**

- Huge data size: >40M documents
- Complex format: lots of digits, tables, formula
- Complex questions: requiring precise knowledge to answer
- Digital sensitive

| | |
|---|---|
| Question | What is the diameter of the **fiber core** in the optical module fiber of HiSecEngine USG6603F-C Tianguan product in international standard? |
| Answer | There are international standard for the diameter of the **fiber core** in optical fibers: **62.5um and 50um for multimode, 9um for single-mode**. |
| Fine-tuned Pangu ❌ | The international standard of the diameter of the **fiber core** in the optical module fiber of HiSecEngine USG6603F-C Tianguan product is **125 micrometers**. |
| +Gauss VectorDB ✅ | The international standard of the diameter of the **fiber core** in the optical module fiber of HiSecEngine USG6603F-C Tianguan product is **62.5um and 50um for multimode, and 9um for single-mode**. |

- Telecom Knowledge QA scores increased **from 10.5 to 96**, reaching the expert level.
- The Pangu model has passed **HCIE** (Huawei Certified ICT Experts), surpassing more than **95%** human experts.



Query → Query Understanding → Query Embedding / Query Keywords → Gauss Vector DB / Keywords Indexing → Answer Composition → Answer

Pangu LLM

HUAWEI  NOAH'S ARK LAB

# Pangu Agent for Operation Research Problem Solver

## Application Case: Supply Chain Optimization

- **Pangu Language Models: Problem Modeling**
  - Dialog
  - Complex reasoning
  - Code generation
  - Summarization

- **OPTV AI Solver (Plug-ins): Problem Solving**
  - Interpreting
  - Calculation
  - Optimization
  - Math problem solving

For level-1 problems of general operation research, we reaches a comparable level with GPT-4.

| Training data | 32,108 |
| --- | --- |
| Test data | 289 |
| GPT-4 | 0.7048 |
| Pangu Production Model | 0.6989 |

The supply chain planner needs to make decisions on the allocation of key components (raw materials), SC (structural component) 1 and 2. There are three requirement codes, namely INV (inverter) A, B, and C. The requirements are 50 pieces, 100 pieces, and 150 pieces, respectively. The virtual material cost is 1.3, 2.1 and 1.8 respectively. The inventory of raw material SC 1 and 2 is 100 pieces and 90 pieces, respectively. SC 1 can be used to produce INV A or B, and processing unit INV A or B requires two SC 1; Processing unit INV C requires one SC 2, and SC 2 can replace SC 1, but SC 1 cannot replace SC 2. May I ask how to make a decision on the production quantity of INVs A, B, and C, as well as the allocation of raw material SC 1 and 2, in order to maximize the total virtual material cost of the requirements met? Build a mathematical planning model for this problem.

Variables:
(Prod. quant. of INV A) x
(Prod. quant. of INV B) y
(Prod. quant. of INV C) z
(The quant. of SC 2 replacing SC 1) r
Objective function:
(Maximizing total virtu. mater. cost) 1.3x+2.1y+1.8z
Constraints:
(Max Prod. quant. constr. of INV A) x<=50
(max Prod. quant. constr. of INV B) y<=100
(max Prod. quant. constr. of INV C) z<=150
(consum. quant. of SC 1 cannot exceed invent.) 2x+2y r<=100
(consum. quant. of SC 2 cannot exceed invent.) z+r<=90
(non neg. int. deci. var. constr.) x, y, z, r \in \mathbb{N})

The solution:
Solution status: OPTIMAL
Target value: 267.0
Decision variables:

| Variable Name | Solution results |
| --- | --- |
| z | 90.0 |
| x | 0.0 |
| y | 50.0 |
| r | 0.0 |

HUAWEI    NOAH'S ARK LAB

# Pangu AI4Industries: 1000+AI Projects in 10+Industries implemented

## Industrial Simulation: COMAC "Oriental Wing Wind"

Aerodynamic AI simulation of large passenger aircraft wing.

The **average error** of flow field prediction is as low as **0.001**.

The single **simulation speed** is improved by **1000 times**.

**SAIL Award** - WAIC2023 Top Award

## Meteorology: Publication in "Nature"

24-hour global weather forecast in **1.4 seconds**.

**10,000 times** faster than traditional methods.

The error is reduced by more than **12%**.

Computing power consumption reduced by **600,000 times**.

Extreme weather forecasts increased by **25%**

## Industrial manufacturing: automatic production scheduling

The time for production line allocation schedule is

reduced **from several hours to 1 minute**.

## EDA

Large Language Model Code Generation.

Test sample generation coverage reaches **99.5%**.

E2E efficiency of test R&D improved by **>3x**

## Drug Discovery

Significantly shorten the drug development cycle.

New broad-spectrum antibiotics were discovered in

Xijiao University Affiliated Hospital **within a month**.

## Government Affairs

Automatic scheduling of thousands of

back-end applications.

Quick realization of various services in cities.

HUAWEI   NOAH'S ARK LAB

# Thank you!

把数字世界带入每个人、每个家庭、
每个组织，构建万物互联的智能世界。

Bring digital to every person, home and organization
for a fully connected, intelligent world.