# EFFICIENTLLM

## Technique

### Architecture Efficiency
- Efficient attention
- Efficient positional encoding
- Sparse modeling via MoEs
- Attention-free sequence modeling

### Data Efficiency
- Data quality and filtering
- Curriculum learning
- Data augmentation

### Inference Efficiency
- Model compression
- Algorithm-Level optimizations
- System-Level optimizations

### Training and Tuning Efficiency
- Scalable training
- Parameter-Efficient Fine-Tuning

### Budget Efficiency
- Scaling behavior and power laws
- Compute-optimal model scaling

## Data

- Dialogue
- Reasoning
- Multilingual
- Vedio
- Image
- Question-Answering
- Text Generation
- Commands/API

## Model

- LLaMA 3 Series
- DeepSeek-R1
- Qwen 2.5 Series
- Phi Series
- Yi
- Mistral
- Stable Diffusion 3.5
- Wan 2.1
- Intern-VL-3
- Qwen-VL 2.5
- LLaVA 1.5
- QvQ-72B

## Efficiency Assessment

| AMU | PCU | AL | AT | AEC | MCR |
|-----|-----|-----|-----|-----|-----|
| Average Memory Utilization | Peak Compute Utilization | Average Latency | Average Throughput | Average Energy Consumption | Model Compression Rate |