

Aurora AI: A Technical Whitepaper on Architecture, Training Pipeline, and Safety Framework

Abstract

Aurora AI is a large-scale, general-purpose language model designed for safe, high-quality reasoning, code generation, and domain-specific assistance. This whitepaper describes the reference design for Aurora AI, including its Transformer-based architecture, core hyperparameters, data curation pipeline, multi-stage training process (pre-training, supervised fine-tuning, and reinforcement learning from human and AI feedback), and safety/alignment mechanisms. We also discuss evaluation methodology, deployment considerations, and known limitations to guide future iterations.

Note: All numerical values (parameter counts, token volumes, etc.) are configurable. This document presents a *plausible* configuration that can serve as a canonical spec for “Aurora AI v1”.

1. Introduction

Recent advances in large language models (LLMs) have enabled systems capable of complex reasoning, multi-step planning, and domain-specific workflows. Aurora AI aims to:

1. Provide **strong general-purpose performance** across natural language, code, and structured data.
2. Support **tool-augmented workflows** (APIs, external search, databases, etc.).
3. Maintain a **safety-first alignment posture**, with explicit guardrails and feedback-driven refinement.

Aurora AI is designed as a decoder-only Transformer model (GPT-style architecture) with a multi-stage training pipeline and modular safety stack, making it suitable for research, enterprise, and productized assistant use cases.

2. System Overview

Aurora AI's design can be summarized as:

- **Model family:** Aurora-Base, Aurora-Instruct, Aurora-RL
 - **Core architecture:** Decoder-only Transformer with grouped query attention, RoPE, and SwiGLU feed-forwards
 - **Scale:** ~34B parameters (reference configuration)
 - **Context window:** 16k tokens (with roadmap to 64k via attention sparsity and compressed memory)
 - **Modalities:**
 - v1: Text + code
 - Optional: Structured tool calls (JSON) and embeddings for retrieval
 - **Deployment:** Sharded inference across GPU or TPU clusters; quantized variants (8-bit / 4-bit) for edge and low-latency serving.
-

3. Model Architecture

3.1 Key Hyperparameters

Component	Value (Aurora-Base-34B)
Parameters	~34 billion
Layers (depth)	64 Transformer blocks
Hidden size	8192
Attention heads	64 (grouped query attention)
FFN expansion (SwiGLU)	4× hidden (\approx 32,768 units)
Vocabulary size	64k BPE tokens
Context length	16,384 tokens
Positional encoding	Rotary (RoPE)

Norm type	Pre-Layer RMSNorm
Activation	SwiGLU
Precision	bf16 (training), fp16/bf16/8-bit (inference)

3.2 Tokenization and Input Representation

- **Tokenizer:** Learned BPE tokenizer with ~64k vocabulary over:
 - English and high-resource languages
 - Programming languages (Python, JS/TS, Go, Rust, etc.)
 - Common markup formats (JSON, YAML, Markdown)
- **Special tokens:**
 - `<bos>`, `<eos>`, `<pad>`, `<unk>`
 - `<tool_call>`, `<tool_result>`, `<sys>`, `<user>`, `<assistant>` for conversation and tool orchestration
- **Segmentation:** The model is trained on mixed single-turn and multi-turn examples, with explicit role tokens to improve dialogue coherence.

3.3 Transformer Block Design

Each layer consists of:

1. **RMSNorm**
2. **Multi-Head Grouped Query Attention (GQA):**
 - K/V heads < Q heads for better memory efficiency
 - Rotary positional embeddings applied in Q/K space
 - Causal attention mask for autoregressive decoding
3. **Residual connection**

4. Second RMSNorm

5. SwiGLU Feed-Forward Network:

- $\text{FFN}(x) = W_3 \left((W_1 x) \odot \sigma(W_2 x) \right)$
- Provides better expressivity than ReLU/GELU at similar compute cost

6. Residual connection

This standardizes the block as:

```
x ← x + MHA(RMSNorm(x))  
x ← x + FFN(RMSNorm(x))
```

3.4 Optimization and Training Infrastructure

- **Optimizer:** AdamW
 - $\beta_1 = 0.9$, $\beta_2 = 0.95$, weight decay = 0.1
- **Schedule:** Cosine decay with warmup
 - Warmup: 2k – 8k steps (configurable)
- **Batching:**
 - Global batch size: ~2–4M tokens/step (data-parallel)
 - Gradient accumulation for effective large batch sizes
- **Parallelism:**
 - Tensor parallelism for large layers
 - Pipeline parallelism across layer partitions
 - Data parallelism across nodes
- **Precision:**

- Computation in bf16
 - Master weights in fp32
 - Gradients clipped to control instability
-

4. Data and Pre-Processing Pipeline

4.1 Data Sources

Aurora AI's pre-training corpus is a mixture of:

- **Public web text** (filtered for quality, deduped)
- **Code corpora** (open-source repositories under permissive licenses)
- **Technical documentation and manuals**
- **Books and long-form articles** (where licensing permits)
- **Dialogue and conversational datasets** (synthetic and human-generated)

A reference configuration targets:

- ~3 trillion tokens total
 - 60–65% general web & natural language
 - 15–20% code and technical artifacts
 - 10–15% instructional & Q/A
 - 5–10% conversational and chat data

4.2 Data Cleaning and Filtering

The pre-processing pipeline includes:

- **Deduplication:** MinHash / LSH-based near-duplicate detection at document and paragraph level.
- **Language identification:** Filtering for primary training languages; low-confidence segments are excluded.
- **Toxicity and policy filtering:**
 - Automatic classifiers (toxicity, hate, sexual content, explicit violence, etc.)
 - Heuristic filters and allow/deny lists.
- **PII minimization:**
 - Heuristic and ML-based PII detectors for names, emails, phone numbers, addresses, etc.
 - Removal or anonymization where feasible.
- **Formatting normalization:**
 - Preservation of code indentation and block structure.
 - Markdown, HTML, and LaTeX cleaning.
 - Collapsing redundant whitespace and boilerplate.

4.3 Mixture Weighting and Curriculum

- Separate **mixtures** are maintained (code, doc, chat, general web).
 - Curriculum schedules:
 - Early stage: higher proportion of general web and books (language broadening).
 - Mid-to-late stages: increased code and instruction density to improve reasoning and tooling.
-

5. Training Regimen

Aurora AI's training is organized into three major phases:

1. **Phase I: Pre-training (Aurora-Base)**
2. **Phase II: Supervised Fine-Tuning (Aurora-Instruct)**
3. **Phase III: Reinforcement Learning (Aurora-RL)**

5.1 Phase I – Pre-Training (Aurora-Base)

Objective: Learn general language and world knowledge through causal language modeling.

- **Objective:** Minimize cross-entropy loss on next-token prediction.
- **Data:** Mixed unlabelled corpus, ~3T tokens.
- **Sequence lengths:**
 - Start with shorter sequences (2k–4k tokens) for stability.
 - Gradually increase to 16k to teach long-context handling.
- **Stabilization techniques:**
 - Warmup + cosine LR schedule.
 - Gradient clipping.
 - EMA tracking and checkpointing.

5.2 Phase II – Supervised Fine-Tuning (SFT) → Aurora-Instruct

Objective: Align the base model with user-centric instructions and multi-turn dialogue patterns.

- **Data sources:**
 - Human-written instruction–response pairs.
 - High-quality curated Q/A sets (math, code, reasoning, explanation).
 - Synthetic data from older Aurora checkpoints filtered by quality and diversity.

- **Task types:**
 - Instruction following and step-by-step reasoning.
 - Code generation and code explanation.
 - Data transformation (JSON, tables, Markdown).
 - Tool calling (API schemas, function call formats).
- **Training recipe:**
 - Smaller LR than pre-training, with limited epochs (to avoid catastrophic forgetting).
 - Dialogue serialization using <sys>/<user>/<assistant> role tokens.
 - Additional loss on tool-call JSON structures to improve reliability.

Result: **Aurora-Instruct** models optimized for helpfulness and structured outputs.

5.3 Phase III – Reinforcement Learning (RLHF / RLAIF) → Aurora-RL

Objective: Align the model with human preferences on safety, helpfulness, and style.

1. **Data Collection:**
 - Human labelers compare pairs of model outputs (A/B) for the same prompt.
 - Labeling criteria: helpfulness, correctness, safety, clarity, style.
 - For low-risk domains, preference data may be generated via *AI feedback* (RLAIF), then spot-checked.
2. **Reward Model (RM):**
 - Separate Transformer model trained to predict human preference scores.
 - Inputs: prompt + candidate output.
 - Outputs: scalar reward.

3. RL Optimization:

- PPO-like objective:
 - Encourage high-reward responses.
 - Penalize large deviations from the SFT policy (via KL penalty).
- Multi-objective reward:
 - $R_{\text{total}} = \alpha * R_{\text{helpfulness}} + \beta * R_{\text{safety}} + \gamma * R_{\text{factuality}}$

4. Safety-specific Alignment:

- Additional training on “rejection” or redirection responses for disallowed content.
- Constitutional-style self-critique where the model revises unsafe drafts.

Result: **Aurora-RL**, Aurora AI’s primary deployment variant with improved safety and preference alignment.

6. Safety, Policy, and Guardrails

Safety is integrated at multiple levels:

6.1 Policy Grounding

- Aurora AI is trained and fine-tuned with reference to a safety policy describing:
 - Disallowed content (e.g., self-harm facilitation, explicit violence, etc.).
 - Sensitive contexts (politics, health, minors).
 - Requirements to avoid hallucinating dangerous instructions or private data.

6.2 Model-Level Safeguards

- RLHF/RЛАIF reward components penalize:

- Policy violations.
- Explicit unsafe instructions.
- Overconfident statements when the model is uncertain.
- SFT includes:
 - Canonical refusal patterns.
 - Explanatory redirections to safer alternatives.

6.3 System-Level Guardrails

At inference time:

- **Input classifiers** to detect high-risk queries.
- **Output filters:**
 - Toxicity and PII detection.
 - Pattern-based filters for known high-risk strings.
- **Tool sandboxing:**
 - Rate limits and scopes on external APIs.
 - Access control for user-specific data.

7. Evaluation and Benchmarking

Aurora AI is evaluated along four major axes:

1. General Language Understanding & Reasoning

- Standard LLM benchmarks (e.g., reading comprehension, QA, logical reasoning tasks).

2. Code and Tool Use

- Code generation and completion tasks.
- Tool-calling correctness: JSON schema compliance, argument accuracy.

3. Safety and Alignment

- Red-teaming on unsafe prompt sets.
- Safety benchmark suites for disallowed content.

4. Human Evaluation

- Side-by-side preference ratings vs. baseline models.
- Task-specific evaluations (e.g., documentation assistance, data analysis workflows).

Typical metrics:

- Pass@k for coding challenges.
 - Win-rate vs baseline in pairwise human eval.
 - Safety violation rate on curated adversarial prompts.
 - Output structure validity (e.g., % of valid JSON tool calls).
-

8. Deployment and Inference

8.1 Serving Architecture

- **Inference engine:** Optimized runtime supporting:
 - Tensor and pipeline parallelism.
 - KV-cache reuse and paged attention.

- **Latency management:**
 - Request batching.
 - Speculative decoding (small draft model + large verifier).
 - Quantization (8-bit or 4-bit weights) for latency-sensitive applications.

8.2 Tooling and Integration

Aurora AI exposes:

- **Chat completion API** for multi-turn conversations.
 - **Tool/function-calling API** using structured JSON descriptions of tools.
 - **Embedding API** (optional Aurora-Embed model) for retrieval-augmented generation (RAG).
-

9. Limitations and Future Work

9.1 Known Limitations

- **Hallucinations:** While reduced by alignment and RAG, Aurora AI can still produce confident but incorrect responses.
- **Biases:** Inherited from training data; ongoing work is required to mitigate harmful stereotypes or unfair associations.
- **Long-context reasoning:** Although 16k tokens are supported, true deep reasoning over extremely long documents remains a challenge.
- **Tool reliance:** Misconfigured external tools can degrade output quality despite correct model behavior.

9.2 Future Directions

- **Extended context windows** ($\geq 64k$) with memory-efficient attention.

- **Stronger retrieval integration** to ground answers in verifiable sources.
 - **Multi-modal extensions** (images, audio, video) with shared embeddings.
 - **More granular safety controls** and user-tunable personas within policy bounds.
 - **Continual learning** loops that safely incorporate user feedback without catastrophic forgetting.
-

10. Conclusion

Aurora AI is designed as a scalable, safety-aware large language model capable of powering a broad range of applications—from general assistants to specialized domain agents. Its architecture and training pipeline follow state-of-the-art practices: large-scale pre-training, instruction tuning, and reinforcement learning from human and AI feedback, all wrapped in a multi-layer safety framework.

This whitepaper provides a reference blueprint for Aurora AI’s parameters, training process, and architecture. It can be adapted to specific deployments (smaller or larger parameter counts, narrower domains, extended modalities) while preserving the core principles of performance, alignment, and safety.