# Contents

- **A quick intro to LLMs**

- **Language Gap and Urdu/Regional Challenges**

- **LLMs Open-Source Ecosystem**

- **Building LLMs for Regional Challenges**

- **Local LLM Architectures**

- **Resource Optimization, Open-Source Fine-Tuning**

- **Applications & Use Cases**

- **Implementation Strategy**

# INTRODUCTION AND CONTEXT

# What are LLMs?

## ARTIFICIAL INTELLIGENCE

Techniques that enables computers to mimic human behavior
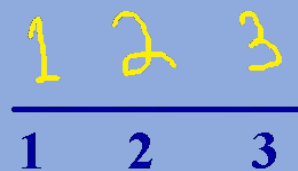
## MACHINE LEARNING

Learning from data to identify patterns and predict with minimal human intervention.

## DEEP LEARNING

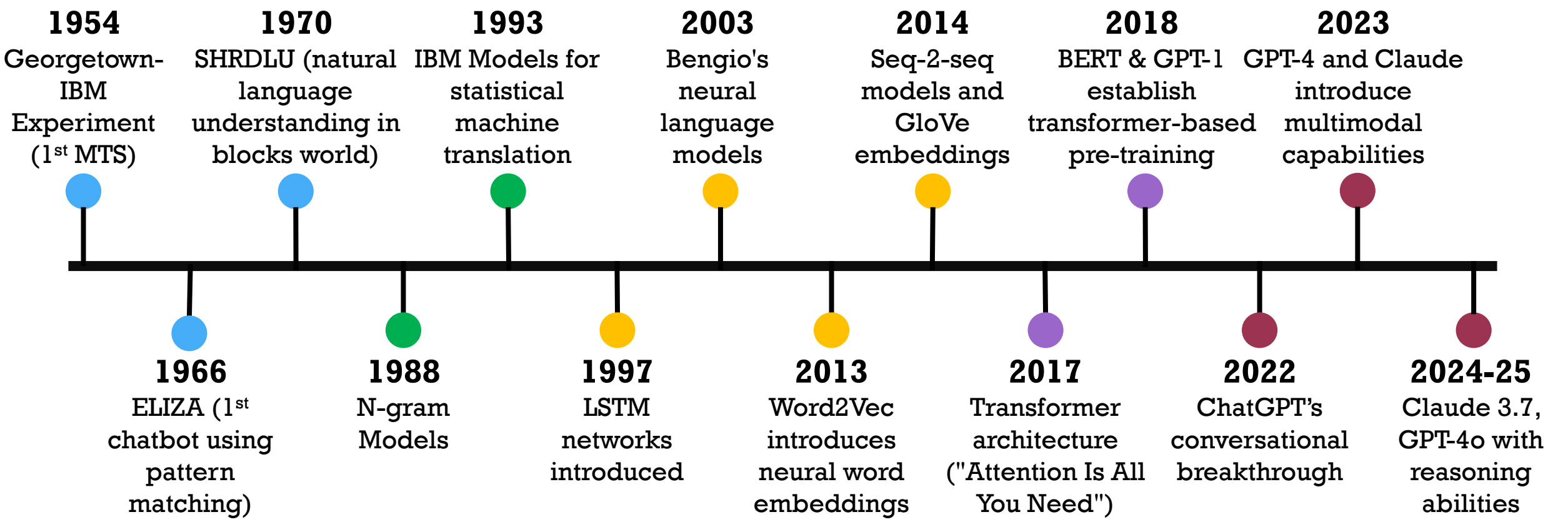Pattern recognition and prediction using neural networks

## LLMs
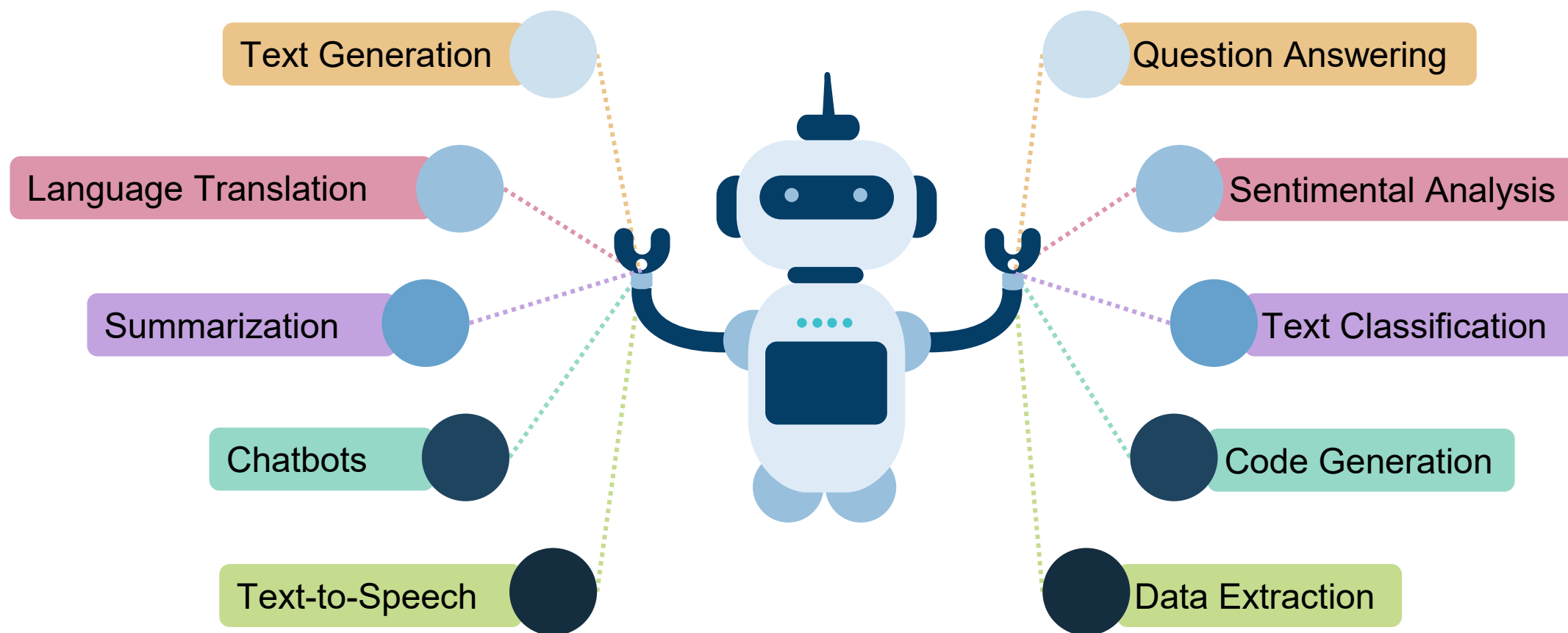
Advanced computer programs that understand and generate human language.

"A Large Language Model is an AI that processes and generates human language using vast text data and deep neural networks.

# What They Can Do?

Text Generation

Language Translation

Summarization

Chatbots

Text-to-Speech

Question Answering

Sentimental Analysis

Text Classification

Code Generation

Data Extraction

# Real-World Impact Across Industries

**All Industries:** LLMs can now automate '60-70% of employees' time. [McKinsey]

**Agriculture:** Decrease operating costs by 22% [ARK]

**Consulting:** 25.1% faster, 40% higher quality [Harvard/BCG]

**Healthcare:** 58.7% better interpersonal skills [Nature]

**Software development:** 55.8% faster [Microsoft/MIT]

**Medicine:** Charts 75% faster, 250% more detail [Carbon Health]

**Finance, accounting, auditing:** Fully exposed, 100% faster [OpenAI/UPenn]

**Therapy/coaching:** AI promises more effective therapy [Seligman]

**Chemistry:** GPT-4 represents the future of chemistry [White]

**Customer service**: 14% more issues resolved [McKinsey]

**Writers:** Fully exposed, 100% faster [OpenAI/UPenn]

❑ **Reasoning:**

  ➢ OpenAI o1/o3: "Thinking time" solves Math Olympiad problems (beating humans)

  ➢ DeepMind's **Gemini Deep Think** hits **gold medal** IMO level

❑ **Multimodal → video next:**

  ➢ **Standard:** Text + vision (GPT-4V, Claude-3, Gemini)

  ➢ **Next:** High-quality video generation (Sora, Veo3, Runway)

❑ **Agentic Use:**

  ➢ Microsoft Copilot: hundreds of millions of Office users

  ➢ Claude used in IDEs like Replit/Cursor for code & automation

❑ **Real-World Tests:** Kaggle AI Chess tournament - OpenAI's o3 beat xAI's Grok-4

❑ **Open source catching up:** LLaMA 3.1, Qwen 2.5, Kimi K2, GLM 4.5 - Beating commercial AIs on key benchmarks
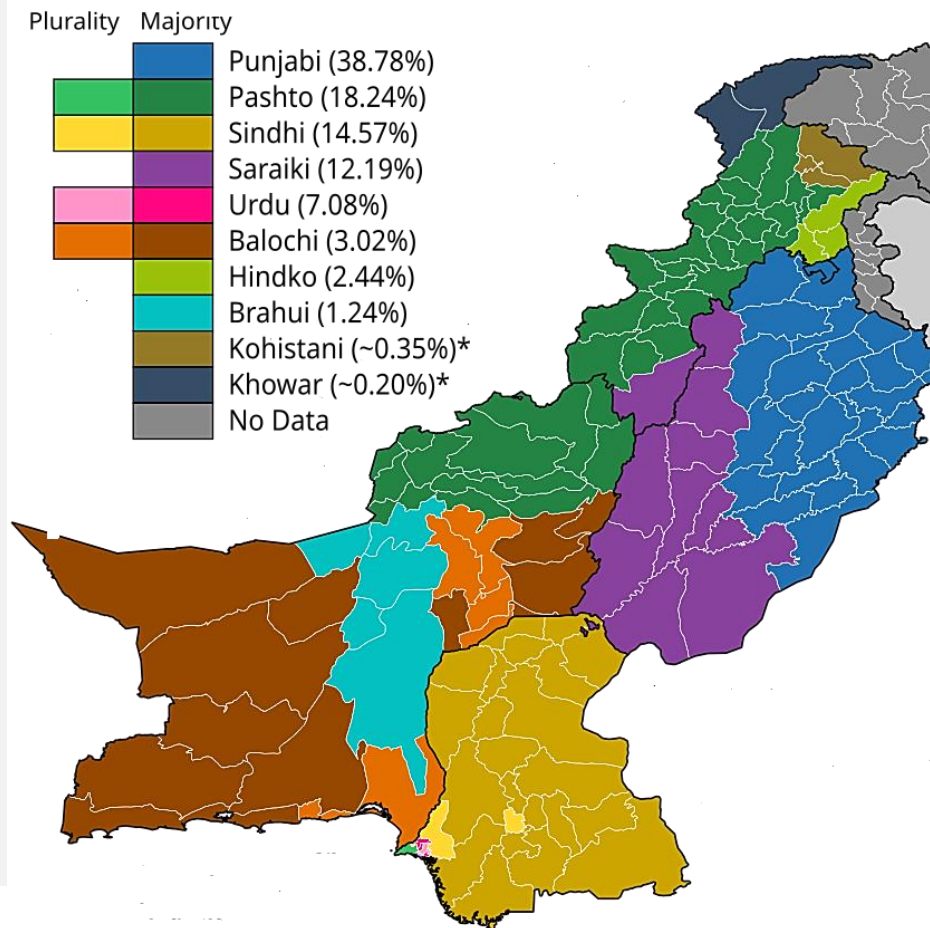
# WHY REGIONAL LLMs?

**The Global Challenge:**

- 7,100+ languages spoken worldwide

- Top 20 languages represent 50% of global speakers

- 98% of AI research focuses on fewer than 10 languages

- English dominates despite being spoken by only 18% of the world population

**The Regional Reality:**

- Millions of people lack access to AI in their native languages

- Critical services, education, and information remain inaccessible

| Plurality | Majority | |
|---|---|---|
| | | Punjabi (38.78%) |
| | | Pashto (18.24%) |
| | | Sindhi (14.57%) |
| | | Saraiki (12.19%) |
| | | Urdu (7.08%) |
| | | Balochi (3.02%) |
| | | Hindko (2.44%) |
| | | Brahui (1.24%) |
| | | Kohistani (~0.35%)* |
| | | Khowar (~0.20%)* |
| | | No Data |

*Pakistan's Linguistic Landscape as of the 2017*
*Source: Wikipedia (Pakistani census 2017)*

**Translation vs. Understanding:**

- Google Translate Urdu accuracy: 67% (compared to 94% for Spanish)

- ChatGPT cultural context accuracy in Urdu: 71%

- Urdu content in major AI training datasets: <0.05%



صبح کے منظر پر ایک نظم لکھیں

Subah ka manzar bahut hi sundar lagta hai,
Suraj dheere dheere asman se nikalta he,
Chidiya log chilane lagte hai,
aur hawa me ek mithas bhar jata hai.

Ask ChatGPT

Urdu

آفتوں کے دور میں
چین کی گھڑی ہے تُو

afton ke dur min
chin ki ghadi hay tu

English

In times of calamity, you are China's watch.

## The Script Challenges

- **Right-to-left complexity:** Arabic script requires 40% more computational resources than Latin scripts
- **Contextual variations:** Single Urdu word can have 12+ different forms depending on position
- **Code-switching reality:** 73% of Pakistani digital users mix Roman-Urdu with English in conversations
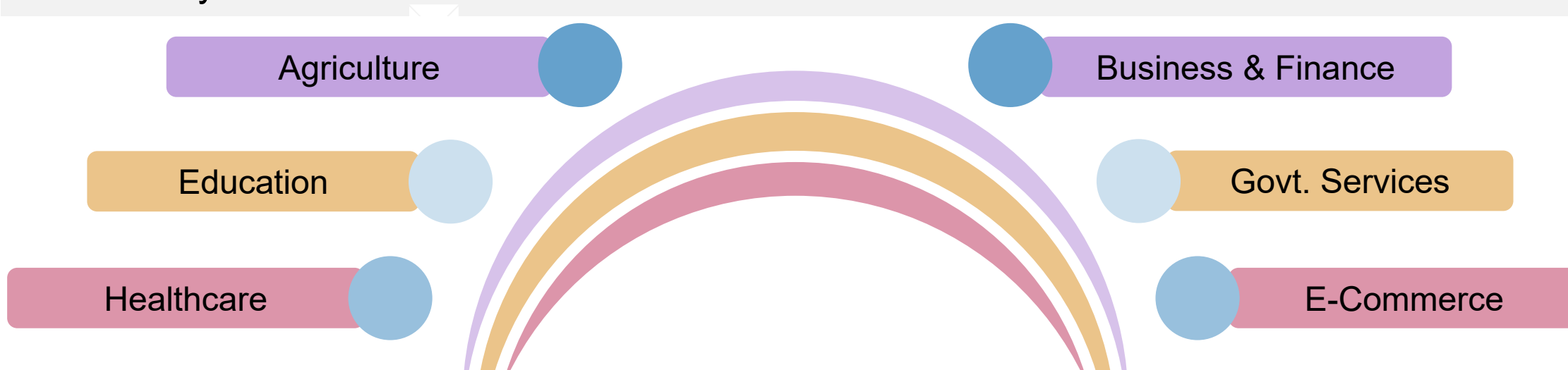
## Data Challenges

- **Training data scarcity:** Urdu represents <0.05% of major AI datasets vs. 54% English
- **Quality gap:** Only 2.3M high-quality Urdu sentences available vs. 570M for Spanish
- **Domain gaps:** Legal, medical, and technical Urdu content virtually non-existent in training data

**What We Need: Pakistan-Native AI Systems**

- **Native Processing:** Fluent in Urdu, Punjabi, Sindhi, and local English patterns.

- **Cultural Intelligence:** Deep understanding of Islamic values, Pakistani traditions, and social norms.

- **Local Expertise:** Knowledge of Pakistani laws, geography, history, and current affairs.

- **Contextual Awareness:** Understanding of regional business practices, educational systems, and daily life.

Agriculture

Business & Finance

Education

Govt. Services

Healthcare

E-Commerce

## Zahanat AI

- **Pakistan's first indigenous GPT** launched by Data Vault (March 2025)

- **Scale**: 1.5B parameters based on Meta's LLaMA architecture

- **Development cost**: <$1M (vs. $5M for DeepSeek)

- **Current status**: Beta testing, limited access within Pakistan

- **Limitations**: English-dominant, basic Urdu support, limited cultural training

## Jazz-NUST-NITB Collaboration (Nov 2024)

- **Scope:** 5-year MoU for indigenous Urdu LLM development

- **Languages:** Primary focus on Urdu, datasets for Pashto and Punjabi

- **Challenge:** Still in early development phase, no public release timeline

**Alif 1.0 by Traversaal AI (Feb 2025)**

- **Performance**: Outperforms Meta LLaMA 3.1 8B on Urdu-specific tasks

- **Innovation**: Native Urdu reasoning with cultural alignment

- **Status**: Open-source model available on HuggingFace

## Still, we have limitations!

- **Market fragmentation:** 6+ separate initiatives with no coordination

- **Limited scale:** Largest model only 8B parameters vs. 175B+ international standards

- **Data scarcity:** <0.1% of internet content in Urdu affects all projects

- **Commercial viability:** No profitable business model demonstrated yet

# Learning from Regional AI Success Stories

## Indic Language Models (India)

- **AI4Bharat's IndicBERT**: Serves 600M users across 22 languages

- **Commercial success**: $12M funding, partnerships with Google, Microsoft

- **Impact**: 340% improvement in local language tasks vs. multilingual models

## Arabic AI Initiatives (Middle East)

- **JAIS Model (UAE)**: $100M investment, 13B parameters trained on Arabic

- **Results**: 89% accuracy in Arabic vs. 34% for GPT-3.5

- **Market capture**: 67% of regional enterprise AI adoption

# THE OPEN-SOURCE ECOSYSTEM

# Commercial vs Open Source LLMs

**Closed Commercial Leaders:**

➤ **OpenAI GPT-4:** $20/month, API access only

➤ **Anthropic Claude:** API-only, usage limits

➤ **Google Gemini:** Integrated products, limited access

**Open-Source Models:**

✓ **DeepSeek V3/R1:** 671B parameters.          30x cheaper than OpenAI o1

✓ **Kimi K2:** 1T parameters.          #1 agentic AI, 90.6% tool success rate

✓ **Qwen 3:** 235B parameters.          119 languages, beats DeepSeek-R1

✓ **GLM 4.5:** 355B parameters.          #3 globally, hybrid reasoning

✓ **OpenAI gpt-oss:**          First open weights since GPT-2 (finally!)

## Open-Source Models are Taking the Lead!

**LLaMA Family (Meta):**

✓ **LLaMA 4** (April 2025): Scout & Maverick variants, 256K context, multimodal

✓ **LLaMA 3.3**: 70B parameters, 128K context window

✓ **Code Llama**: Programming specialist, GitHub integration

**Legacy**: Started the open-source revolution, 40M+ downloads
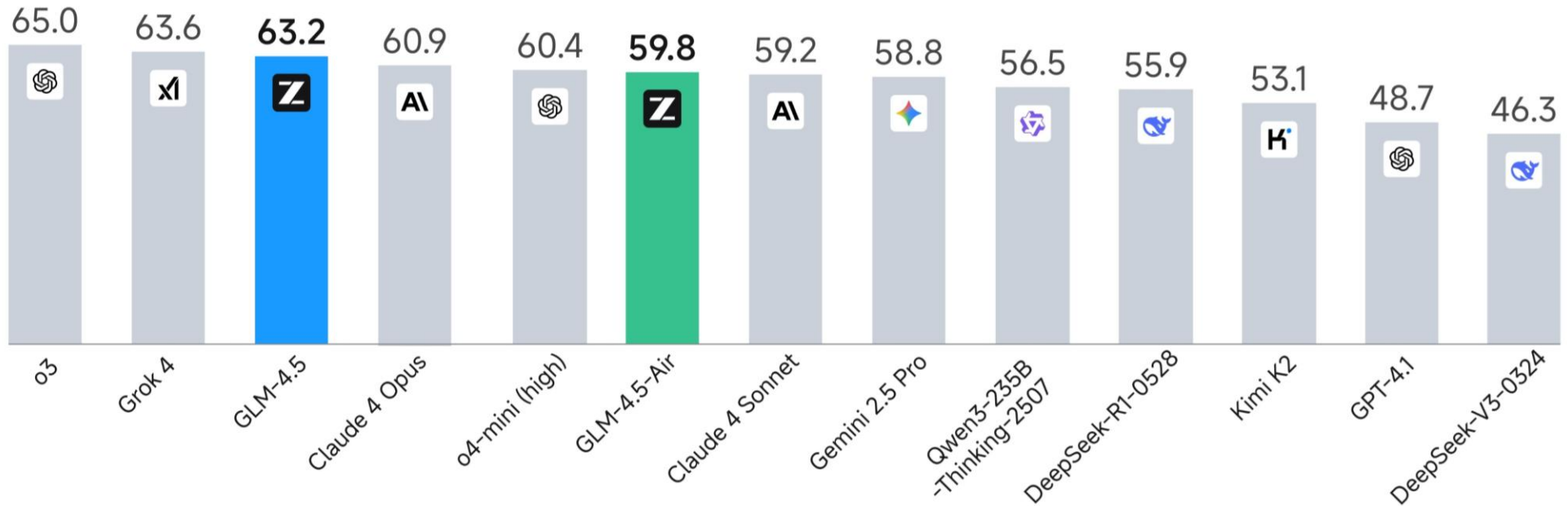
**Mistral Models (Mistral AI):**

✓ **Mixtral 8x22B**: Latest MoE, 141B total parameters

✓ **Mistral Large**: Competing with GPT-4 class models

# Chinese Dominance (2025 Breakthrough)

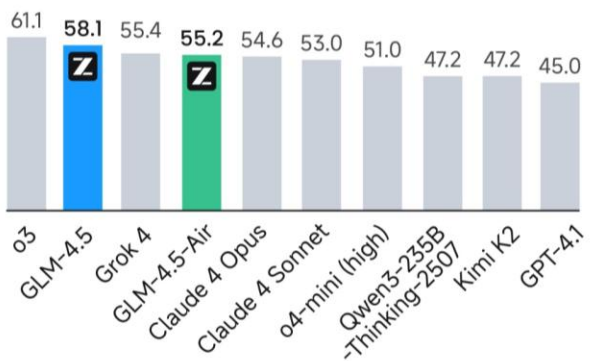| Family | Company | Model | Parameters | Features & Achievements |
|--------|---------|-------|------------|------------------------|
| **Qwen** | Alibaba | Qwen 3 | 235B MoE | 119 languages, Beats DeepSeek-R1 |
|  |  | Qwen 2.5-Max | - | Outperforms DeepSeek V3 in security |
| **DeepSeek** | - | DeepSeek V3 | 671B MoE | Industry-leading cost efficiency |
|  |  | DeepSeek R1 | - | Reasoning specialist, 30x cheaper than o1 |
| **Kimi** | Moonshot AI | Kimi K2 | 1T | Agentic AI leader, 90.6% tool success rate |
| **GLM** | Z.ai | GLM 4.5 | 355B | Hybrid thinking modes, #3 globally ranked |
|  |  | GLM 4.5-Air | 106B | Efficient for edge deployment |

**OpenAI gpt-oss (August 2025):** 120B & 20B variants, Apache 2.0 license

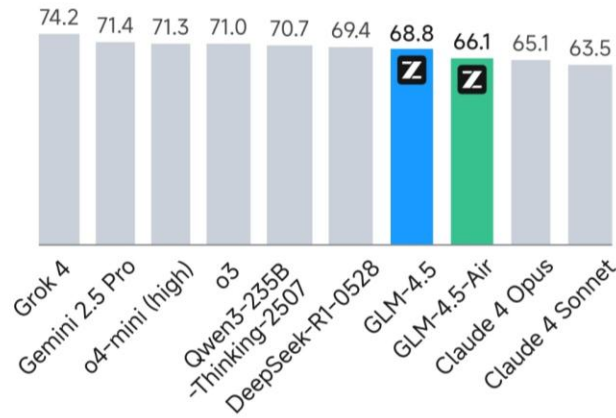**Significance**: Big Tech finally joins open-source movement!

# Chinese Dominance (2025 Breakthrough)



**Main chart (overall scores):**

| Model | Score |
|---|---|
| o3 | 65.0 |
| Grok 4 | 63.6 |
| GLM-4.5 | 63.2 |
| Claude 4 Opus | 60.9 |
| o4-mini (high) | 60.4 |
| GLM-4.5-Air | 59.8 |
| Claude 4 Sonnet | 59.2 |
| Gemini 2.5 Pro | 58.8 |
| Qwen3-235B-Thinking-2507 | 56.5 |
| DeepSeek-R1-0528 | 55.9 |
| Kimi K2 | 53.1 |
| GPT-4.1 | 48.7 |
| DeepSeek-V3-0324 | 46.3 |

## Agentic

Agenic Benchmarks: TAU-Bench, BFCL V3 (Full), BrowseComp

| Model | Score |
|---|---|
| o3 | 61.1 |
| GLM-4.5 | 58.1 |
| Grok 4 | 55.4 |
| GLM-4.5-Air | 55.2 |
| Claude 4 Opus | 54.6 |
| Claude 4 Sonnet | 53.0 |
| o4-mini (high) | 51.0 |
| Qwen3-235B-Thinking-2507 | 47.2 |
| Kimi K2 | 47.2 |
| GPT-4.1 | 45.0 |

## Reasoning

Reasoning Benchmarks: MMLU-Pro, AIME 24, MATH 500, SciCode, GPQA, HLE, LCB (2407-2501)

| Model | Score |
|---|---|
| Grok 4 | 74.2 |
| Gemini 2.5 Pro | 71.4 |
| o4-mini (high) | 71.3 |
| o3 | 71.0 |
| Qwen3-235B-Thinking-2507 | 70.7 |
| DeepSeek-R1-0528 | 69.4 |
| GLM-4.5 | 68.8 |
| GLM-4.5-Air | 66.1 |
| Claude 4 Opus | 65.1 |
| Claude 4 Sonnet | 63.5 |

## Coding

Coding Benchmarks: SWE-Bench Verified, Terminal-Bench

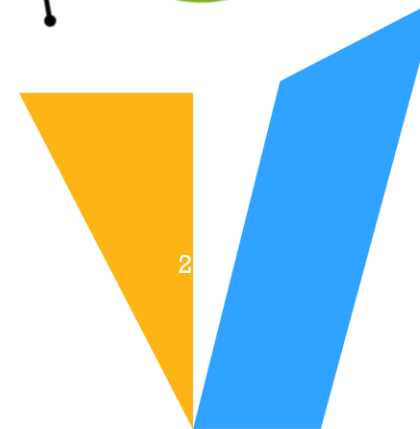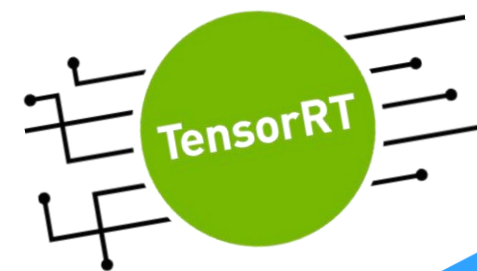| Model | Score |
|---|---|
| Claude 4 Opus | 55.5 |
| Claude 4 Sonnet | 53.0 |
| GLM-4.5 | 50.9 |
| o3 | 49.7 |
| Kimi K2 | 45.2 |
| GLM-4.5-Air | 41.5 |
| GPT-4.1 | 39.5 |
| Gemini 2.5 Pro | 37.2 |
| o4-mini (high) | 36.7 |

**Hugging Face:** 500K+ models, 100K+ datasets, 1M+ developers, 10B+ downloads

**Model Serving & Deployment**

❑ **vLLM:** High-throughput inference server, 24x faster than vanilla PyTorch

❑ **Ollama:** Local model management with one-command deployment

❑ **TensorRT-LLM:** NVIDIA optimized serving for hardware acceleration

❑ **LM Studio:** User-friendly interface for no-code local AI deployment

**Training Frameworks:**

✓ **Transformers (Hugging Face):** Industry standard with 180K+ GitHub stars

✓ **DeepSpeed (Microsoft):** Distributed training for trillion-parameter models

✓ **Accelerate:** Simplified multi-GPU training framework.

# BUILDING REGIONAL MODELS

## Phase 1: R&D

- **Data Collection**: Build comprehensive Urdu corpus
- **Model Development**: Complete initial training experiments
- **Baseline Testing**: Establish performance metrics against existing solutions

## Phase 2: Testing

- **Domain Applications**: Test in education, healthcare, agriculture
- **User Feedback**: Deploy with limited user groups
- **Model Refinement**: Iterate based on real-world performance

## Phase 3: Market Entry

- **Production Deployment**: Launch stable commercial version
- **Partnership Development**: Establish B2B and government relationships
- **Scale Planning**: Prepare for broader market adoption

24

**The Current Reality**

- **Digital Urdu content available**: ~15TB across all sources

- **International LLM standard**: 500TB+ high-quality text required

- **Data collection costs**: PKR 20 per 1K tokens vs. PKR 0.80 for English content

## Strategy

✓ Literary Corpus: Digitize 847 classic and modern Urdu books

✓ News Media: Collect 1.2M articles from 15 major publications (2018-2024)

✓ Social Media: 50M posts with cultural context annotations

✓ Academic Content: Gather 180K research papers from Pakistani universities

✓ Religious sensitivity: Develop Islamic jurisprudence knowledge base

✓ Social understanding: Map family structures and respect hierarchies

✓ Regional diversity: Include Punjabi, Sindhi, Balochi cultural contexts

## Model Selection

## Training

## Benchmarking

- **Open source LLaMA 2-7B Models**
- **Mistral Models**
- **Custom tokenizer**

- **Continual Pre-training**
- **Cultural Fine-tuning**
- **Multi-stage approach**

**Urdu comprehension**: Aim to exceed current GPT-4 performance on Urdu tasks

**Code-switching**: Handle mixed Urdu-English communication patterns

**Three Stage Pipeline**

| Stage | Objective | Outcome |
|---|---|---|
| **Pre-training** | Learn language from internet | General intelligence |
| **Supervised Fine-tuning** | Learn to follow instructions | Conversational ability |
| **RLHF** | Align with human values | Safe, helpful responses |

**The Data Diet**

✓ Urdu News media, web pages, books, wikipedia, literary content, internet and social media etc.

✓ Max no. of tokens of human Urdu text

**Learning Progression**

**Phase 1:** Language structure and syntax

**Phase 2:** Factual knowledge acquisition

**Phase 3:** Reasoning capabilities

**Computational Requirements**

➢ Multi GPU cluster deployment

➢ Continuous training for 3-6 months

**Supervised Fine-tuning (SFT):**

**Goal:** Make model follow specific instructions

**How:** Feed curated human-written local language examples (input → correct output)

Model starts producing task-relevant responses after SFT

**Why SFT matters:**

Pre-training makes a **knowledgeable model**, but not necessarily an **obedient assistant**.

SFT shapes *what the model does with its knowledge*.

**Process:**

**Collect preferences** – Experts rank different model outputs

**Train reward model** – learns what "good" means

**Policy optimization** – model adapts to produce preferred answers
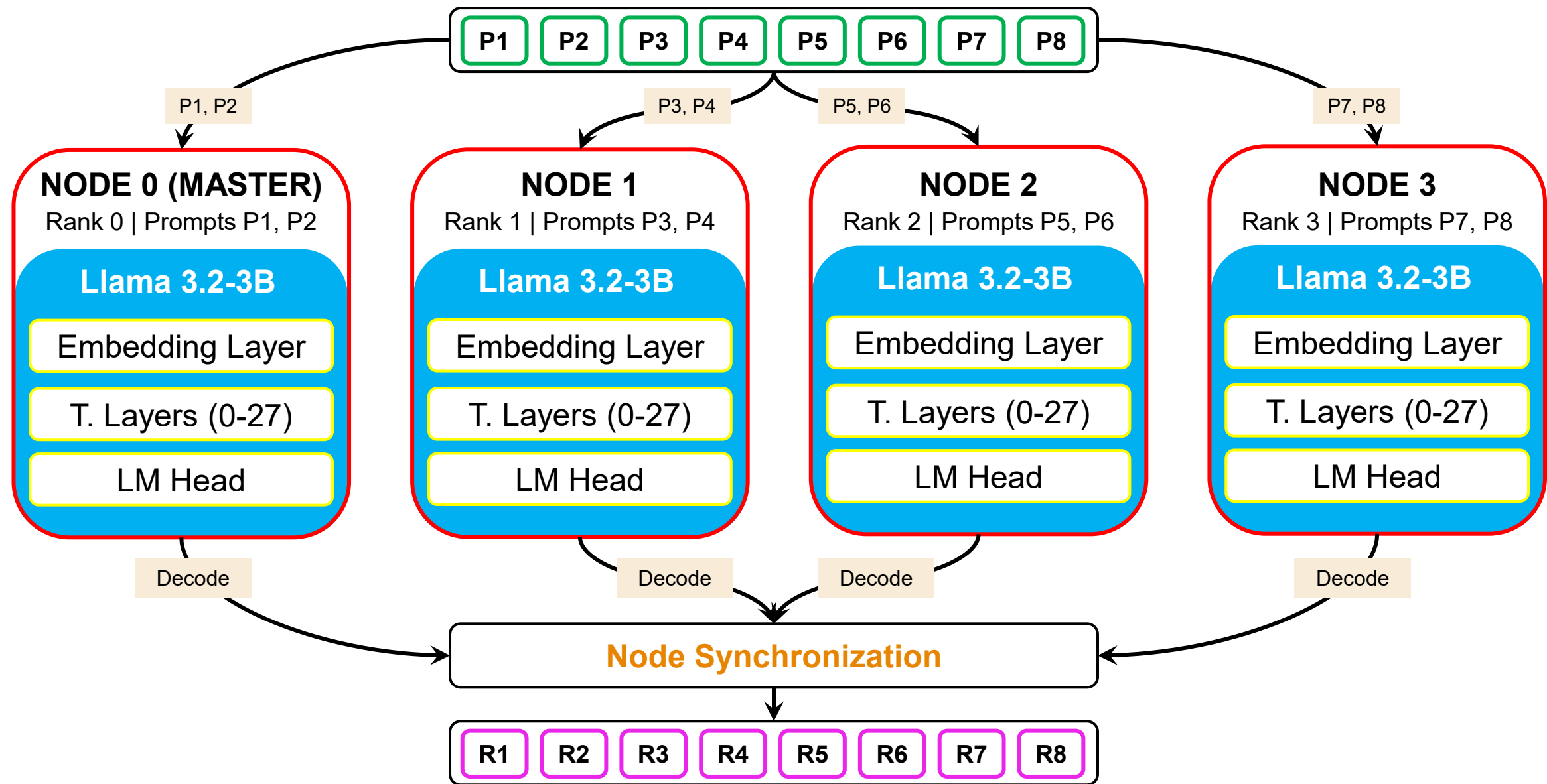
**Outcome:**

- ✓ Safer, more helpful, less biased responses
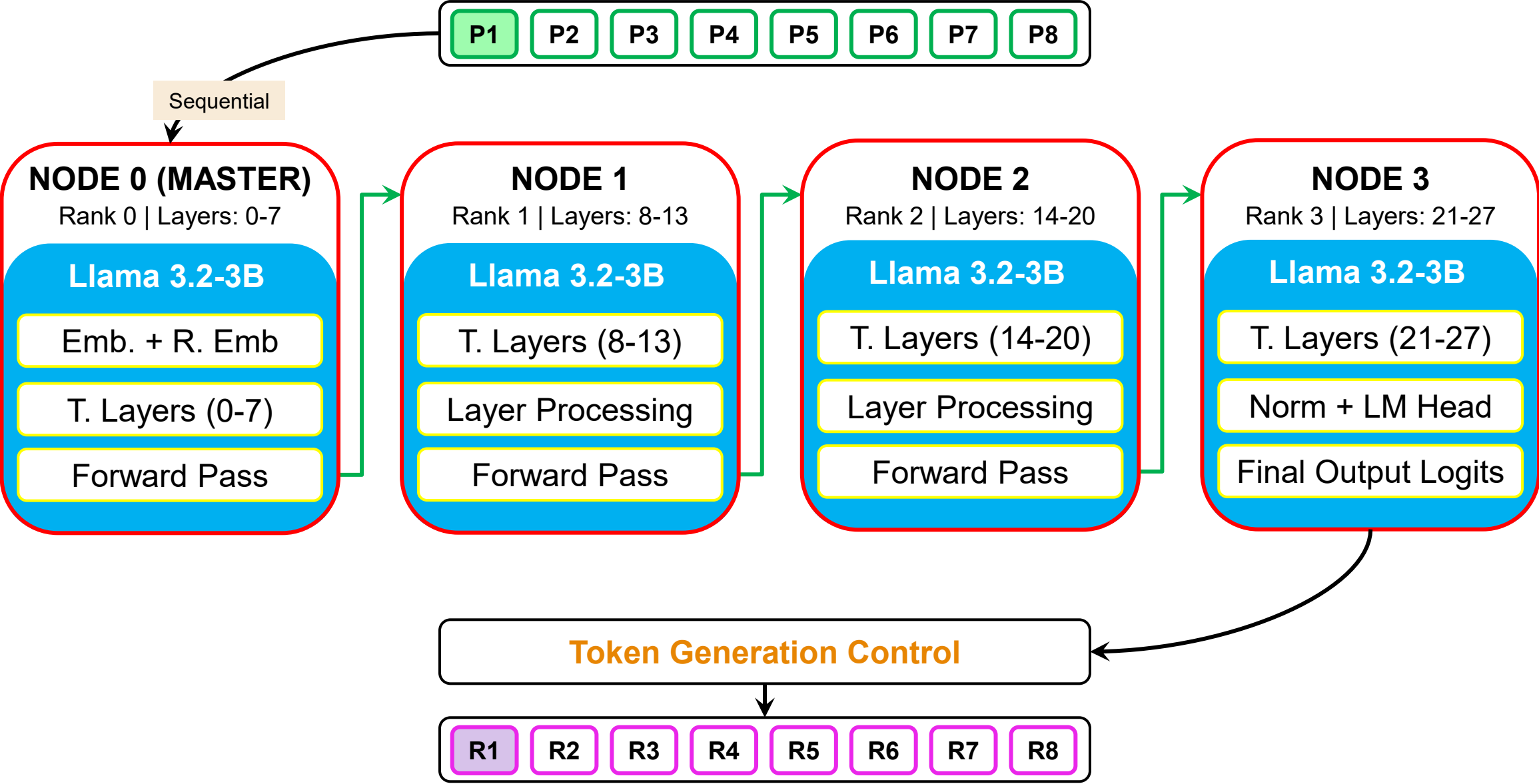
- ✓ Behaves closer to expectations

**Parallelism Strategies:**

- **Data Parallelism:** Replicate model per GPU → train on different batches → sync gradients.

- **Model Parallelism:** Split model layers across GPUs → share parameters via communication.

- **Pipeline Parallelism:** Assign different layers to GPUs → process multiple micro-batches at once.

- **Tensor Parallelism:** Split single ops (e.g., attention) across GPUs → high communication cost.

**Example – Namal HPC Cluster:**

- ✓ 6-node GPU cluster with SLURM job management

- ✓ GPUs: RTX 4070 Ti (12GB VRAM each) = 72GB total GPU memory

- ✓ Successfully implemented distributed training and inference of open source LLMs.

**Model Loading → Tokenization → Layer Processing → Token Generation → Output Decode**

**Cold Start:** NVMe → RAM → GPU (2-5s), CUDA initialization (200-500ms)

**Input Processing:** Tokenization (5-10ms), embedding lookup (1-3ms)

**Core Computation:** Attention $O(n^2)$ operations (50-200ms/layer)

**Token Generation:** Sampling/softmax (2-5ms/token), KV cache updates (1-3ms/token)

**Output:** Text decoding (1-5ms), result aggregation (5-20ms)

**Network Backbone:** InfiniBand (1-2µs) vs Ethernet (µs-ms range)

| Configuration | Primary Overhead | Latency Range | Hardware Cause |
|---|---|---|---|
| **Single Node** | Memory bandwidth | 20-80ms/layer | HBM weight loading |
| **Data Parallel** | Node synchronization | 50-200ms | Coordination barriers |
| **Model Parallel** | Network communication | 5-500ms/layer | Inter-node data transfer |

❑ **Quantization - Reducing Numerical Precision: (**Typically <2% performance degradation**)**

   **FP32 → INT8:** 4x memory reduction, 2-3x speed improvement

   **FP32 → INT4:** 8x memory reduction, significant speed gains

❑ **Pruning - Removing Parameters: (**50-90% parameter reduction maintaining performance**)**

   **Unstructured:** Remove individual weights based on magnitude

   **Structured:** Remove entire neurons, attention heads, or layers

❑ **Distillation - Knowledge Transfer: (**Large "teacher" model trains smaller "student" model**)**

   **Example:** DistilBERT (66M params) achieves 97% of BERT (110M params) performance

**Trade-offs:** All methods balance model size, inference speed, and output quality

# LAUNCHING REGIONAL MODELS

**Market Size Context**

- **Pakistani digital economy**: PKR 2.8 trillion with 34% annual growth

- **Current language gap**: Regional language AI market largely unaddressed

- **Long-term opportunity**: Estimated PKR 580 billion market by 2028

**Potential Revenue Models**

- **Enterprise Solutions**

- **Government Contracts**

- **API Services**

**Education Sector**

**Target market**: 35M Urdu-medium students currently underserved

- ✓ Curriculum-aligned content generation
- ✓ Personalized learning explanations
- ✓ Support for students struggling with English-medium materials

**Agriculture**

**Target users**: 25M farming families needing localized information

- ✓ **Information gaps**: Weather, crop, and market data rarely available in local languages
- ✓ **Economic potential**: Optimized farming advice could reduce input costs

Local LLMs are the way forward!

# Thank You!