

# MULTILINGUAL MCP

Making tool calling work for the next billion users



Samyuktha M S  
Software Developer, IBM

MCP Dev Summit, Bangalore  
10 June 2026

# Why I'm here ?

---



Hi, I'm **Samyuktha**

- Software Developer at IBM, ISL
- AI Enthusiast
- *15x Hackathon winner*
- Notable wins - *Google Agentic AI Hackathon, GrabHack 2.0, Chhalaang 4.0, Google Build and Blog Marathon*
- Speaker at GHCI '25, GHC '26, OpenSearch Con India '26, Devfest
- [linkedin.com/in/samyuktha-m-s](https://www.linkedin.com/in/samyuktha-m-s)
- [github.com/samyuktha-12](https://github.com/samyuktha-12)



# “ನನ್ನ ಮಗನಿಗೆ ಜ್ವರ ಬಂದಿದೆ”

My son has a fever

Meet Lakshmi.

- 34, lives in Mandya district,
- four hours from Bengaluru by bus.
- Speaks Kannada. Only Kannada.

- Last year her panchayat publicised a government health helpline
- **A voice agent.**
- Three tools — diagnose, guide, dispatch. **Powered by MCP under the hood.**



Then a real user - Lakshmi, called it. It failed her.

**Three different ways**

# Break 01: The Wrong Tool

---

00:00.0 Lakshmi: "ನನ್ನ ಮಗನಿಗೆ ಜ್ವರ ಬಂದಿದೆ, ಔಷಧಿ ಬೇಕು"  
My son has fever, I need medicine.

00:00.6 ..... [tool selected: schedule\_appointment]

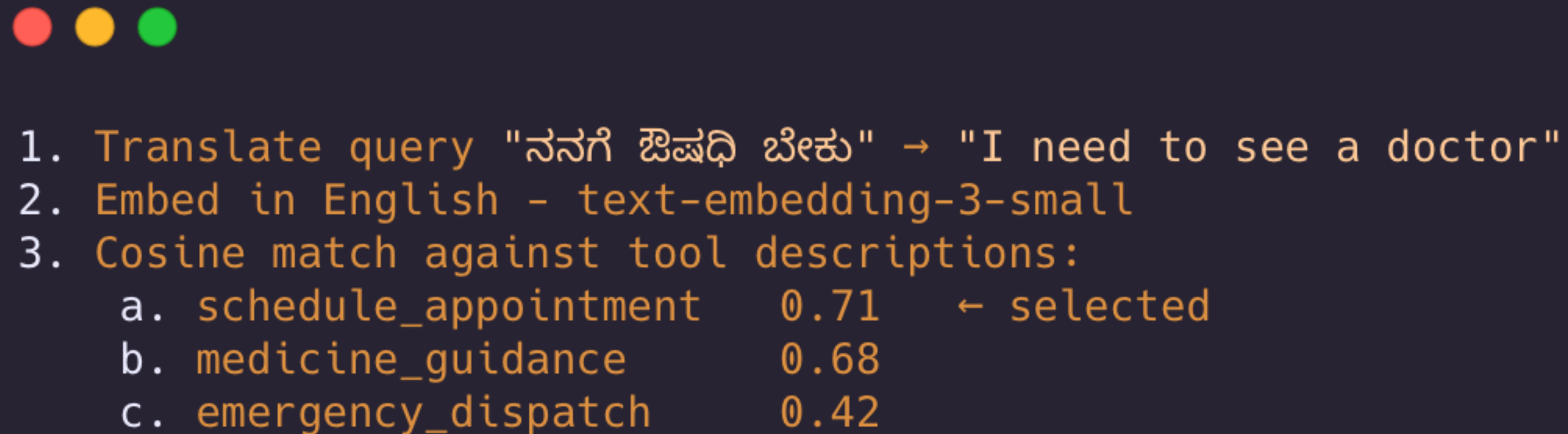
00:00.9 Agent: "ನಿಮ್ಮ ಅಪಾಯಿಂಟ್‌ಮೆಂಟ್ ಮುಂದಿನ ಮಂಗಳವಾರಕ್ಕೆ ಬುಕ್ ಆಗಿದೆ"  
Your appointment is booked for next Tuesday at 3 PM.

She asked for medicine **right now**. She got an appointment for **next week**.  
The **intent** was right. The **tool** was wrong.

## DIAGNOSIS 01

# English-Biased Tool Selection

When the query arrives, three things happen — all in English embedding space.

- 
1. Translate query "ನನಗೆ ಔಷಧಿ ಬೇಕು" → "I need to see a doctor"
  2. Embed in English - text-embedding-3-small
  3. Cosine match against tool descriptions:

a. schedule_appointment	0.71	← selected
b. medicine_guidance	0.68	
c. emergency_dispatch	0.42	

The translation collapsed her urgency.

"ಔಷಧಿ ಬೇಕು" → "see a doctor" is wrong by 0.03 cosine.

In English, both sound reasonable. In Kannada, only one of them is what she said.

# FIX #1 - Multilingual Tool Manifests

Two changes. One in the manifest. One in the retrieval layer.

## Tool manifest

```
{
  "name": "medicine_guidance",
  "descriptions": {
    "en": "OTC medicine recommendations and dosage",
    "kn": "ಔಷಧಿ ಸಲಹೆ ಮತ್ತು ಡೋಸೇಜ್ ಮಾರ್ಗದರ್ಶನ",
    "hi": "दवा सलाह और खुराक मार्गदर्शन",
    "ta": "மருந்து ஆலோசனை மற்றும் அளவீடு வழிகாட்டுதல்"
  }
}
```

## Retrieval layer

```
Index with BGE-M3 (100+ languages, no fine-tune)
Same query "ನನಗೆ ಔಷಧಿ ಬೇಕು" → no translation
Embed directly. Cosine match in shared space.
```

## The Result

- medicine\_guidance cosine 0.68 → 0.91
- Tool selection accuracy 38% → 87%

# Break 02: The Code-Switch

---

00:00.0 Lakshmi: "ನನಗೆ paracetamol ಬೇಕು, ಎಷ್ಟು dose ಕೊಡಬೇಕು bachhe-ge? "  
I need paracetamol – how much should I give my child?

00:00.7 ..... [lang detect: kn 0.51, en 0.32, hi 0.17]

00:00.7 ..... [tokenizer: "para", "##ceta", "##mol" split]

00:00.8 ..... [tool selected: generic\_query]

00:01.3 Agent: "I'm not sure I understood. Could you say that again?"

Kannada grammar. English drug name. Latin script.  
Three writing systems in one sentence — and the system spoke none of them.

## DIAGNOSIS 02

# Language is a token-level property

Standard pipelines treat language as a sentence-level fact: "this sentence is in Kannada." Real sentences don't work that way.

ನನಗೆ  
KN  
pronoun

paracetamol  
EN  
proper-noun

ಬೇಕು  
KN  
verb

ಎಷ್ಟು  
KN  
quant

dose  
EN  
noun

ಕೊಡಬೇಕು  
KN  
verb

bachhe  
HI  
loan

-ge  
KN  
postp.

Sentence-level detection said: 51% Kannada.

But "paracetamol" and "dose" carry the medical meaning.

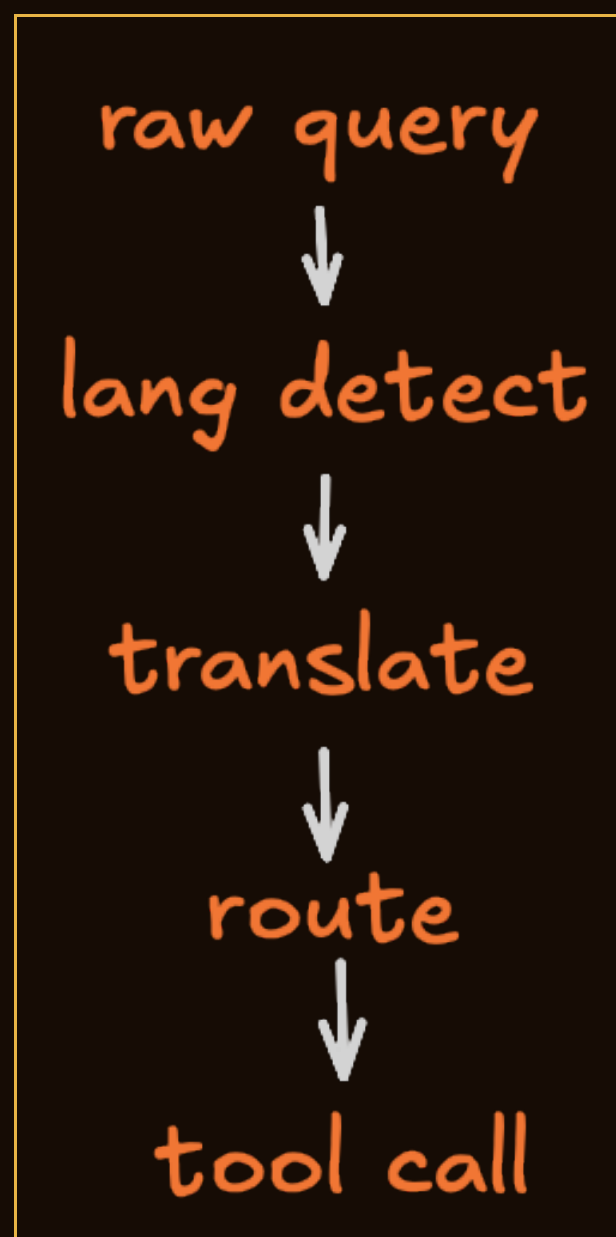
Routing to Kannada loses the English-anchored entities.

Stop detecting language. Detect intent.

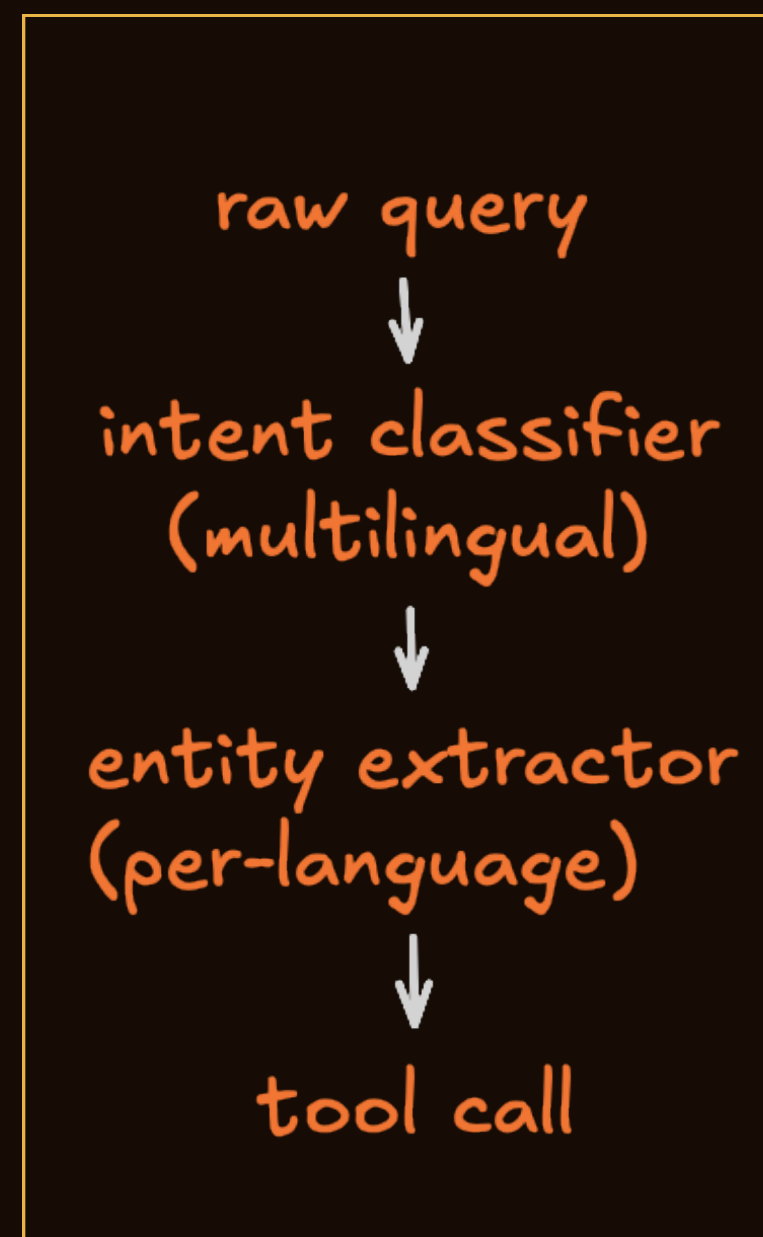
# FIX #2 - Intent-First Routing

Skip the language detection step. Route on intent, preserve entities by language.

Before



After



The Router

```
def route(query: str) -> ToolCall:
    intent = classify(query) # multilingual BGE-M3
    entities = extract(query, per_lang=True)
    return ToolCall(intent.tool, entities)
```

The Result

- Code-switch query accuracy 12% → 84%
- Entity recall 31% → 92%

# Break 03: The English Response

---

```
00:00.0 Lakshmi: "ನನ್ನ ಮಗನಿಗೆ ಜ್ವರ ಬಂದಿದೆ, ಔಷಧಿ ಬೇಕು"  
My son has fever, I need medicine.  
  
00:00.7 ..... [tool: medicine_guidance] ✓  
  
00:00.9 tool_out: {  
    "medicine": "paracetamol",  
    "dosage": "5ml",  
    "frequency": "every 6 hours",  
    "max_per_day": 4,  
    "warning": "consult doctor if fever exceeds 3 days"  
}  
  
00:01.4 Agent: "Give paracetamol, 5 milliliters, every 6 hours, maximum 4 doses per day. Consult  
a doctor if fever exceeds 3 days."
```

The right tool fired. The right data came back. The agent spoke it back in **English**.  
Lakshmi only speaks **Kannada**.

## DIAGNOSIS 03

# Tool Outputs Are Foreign Objects

Look closely at what came back. Four hidden English assumptions in one tool response.

```
{
  "medicine": "paracetamol", ← entity OK
  "dosage": "5ml", ← unit convention
  "frequency": "every 6 hours", ← English phrase
  "max_per_day": 4, ← key in English
  "warning": "consult doctor if fever exceeds 3 days"
}
```

English prose

Some of these can be translated word-for-word. Others can't. "5ml" is English convention. "every 6 hours" is English grammar. The warning string is English prose. The keys are English schema.

Outputs need contracts, not translation.

# FIX #3 - Output Localization Contracts

The tool author declares how each field speaks in each language.

## Manifest Extension

```
{
  "output_schema": {
    "medicine": {
      "verbalize": {
        "en": "{value}",
        "kn": "{value}"
      }
    },
    "dosage": {
      "verbalize": {
        "en": "{value}",
        "kn": "{value_ml_to_spoon}"
      }
    },
    "frequency": {
      "verbalize": {
        "en": "every {hours} hours",
        "kn": "{hours} ಗಂಟೆಗೊಮ್ಮೆ"
      }
    },
    "warning": {
      "verbalize": {
        "en": "{value}",
        "kn": "{value_kn}"
      }
    }
  }
}
```

## Before

"Give paracetamol, 5 milliliters, every 6 hours, maximum 4 doses per day. Consult a doctor if fever exceeds 3 days."

## After

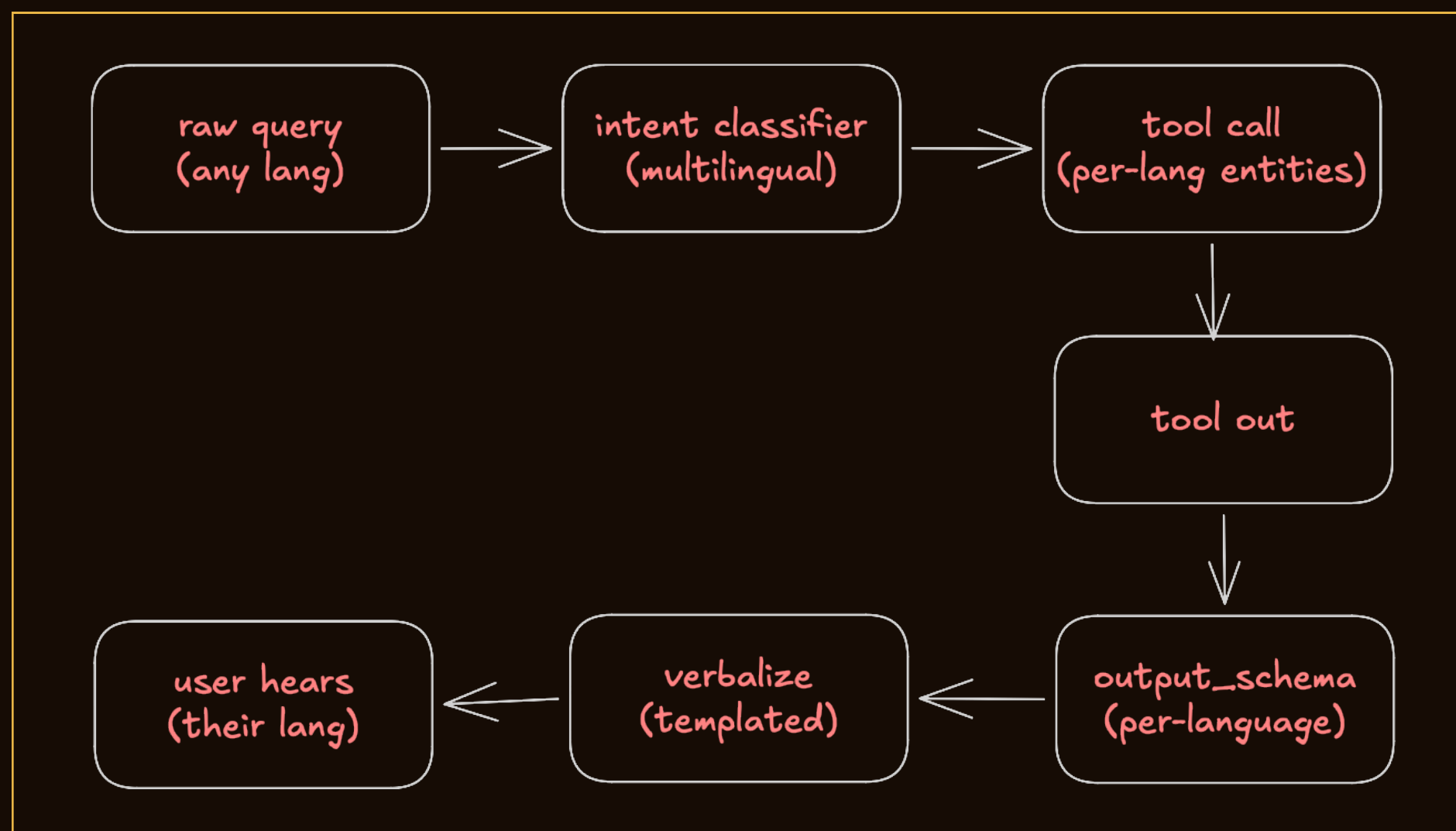
"ಪ್ಯಾರಸಿಟಮಾಲ್ ಕೊಡಿ, ಒಂದು ಚಮಚ, 6 ಗಂಟೆಗೊಮ್ಮೆ, ದಿನಕ್ಕೆ 4 ಬಾರಿ. 3 ದಿನಗಳಲ್ಲಿ ಜ್ವರ ಕಡಿಮೆಯಾಗದಿದ್ದರೆ ವೈದ್ಯರನ್ನು ಸಂಪರ್ಕಿಸಿ."

## The Result

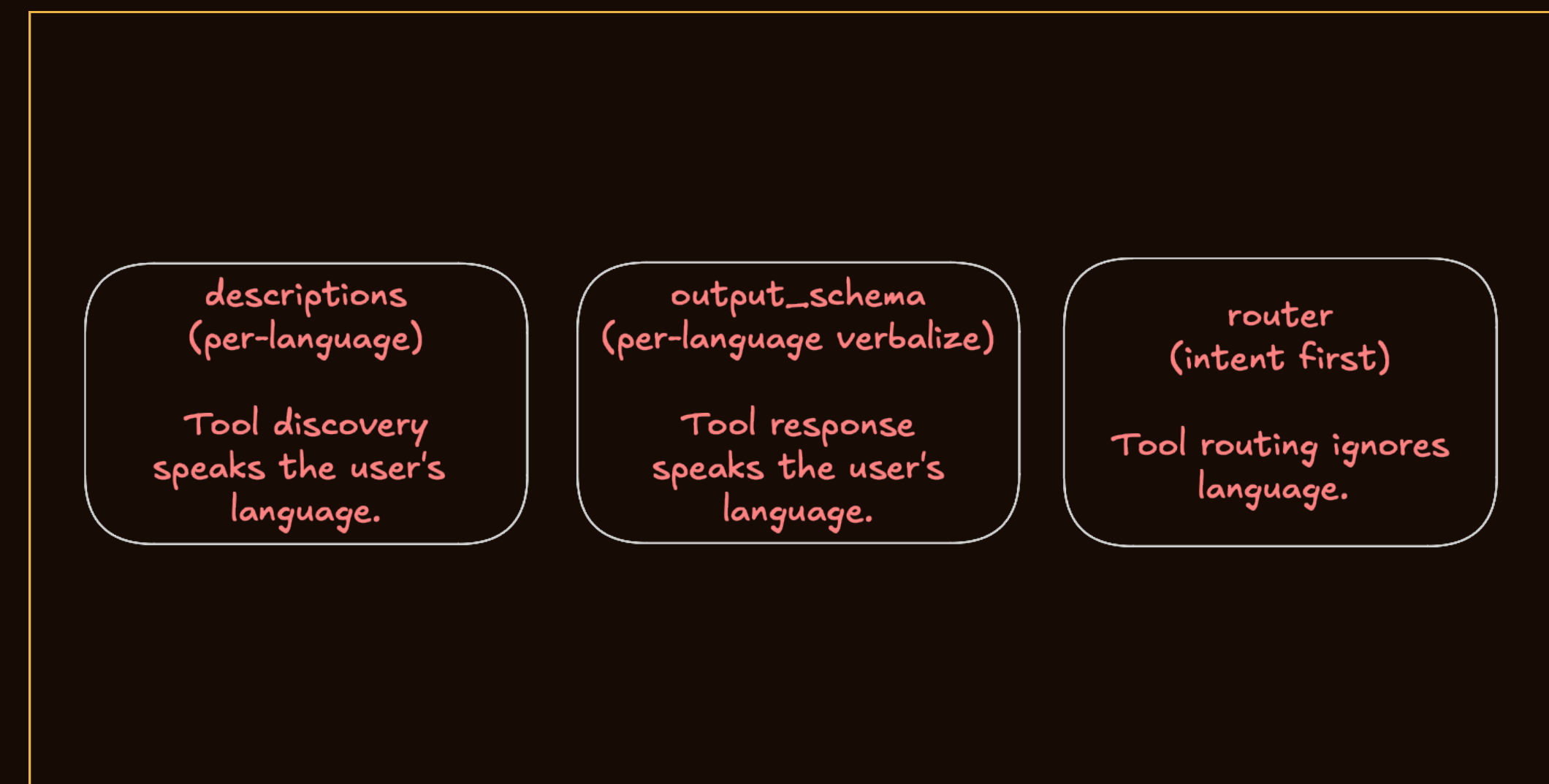
- User completion rate (Kannada) 22% → 71%
- Localized response coverage 0% → 94%

# Lakshmi's full system

## ARCHITECTURE



## THREE MANIFEST EXTENSIONS



Multilingual MCP needs three things the spec doesn't have:  
discovery, routing, response — all in the user's language.

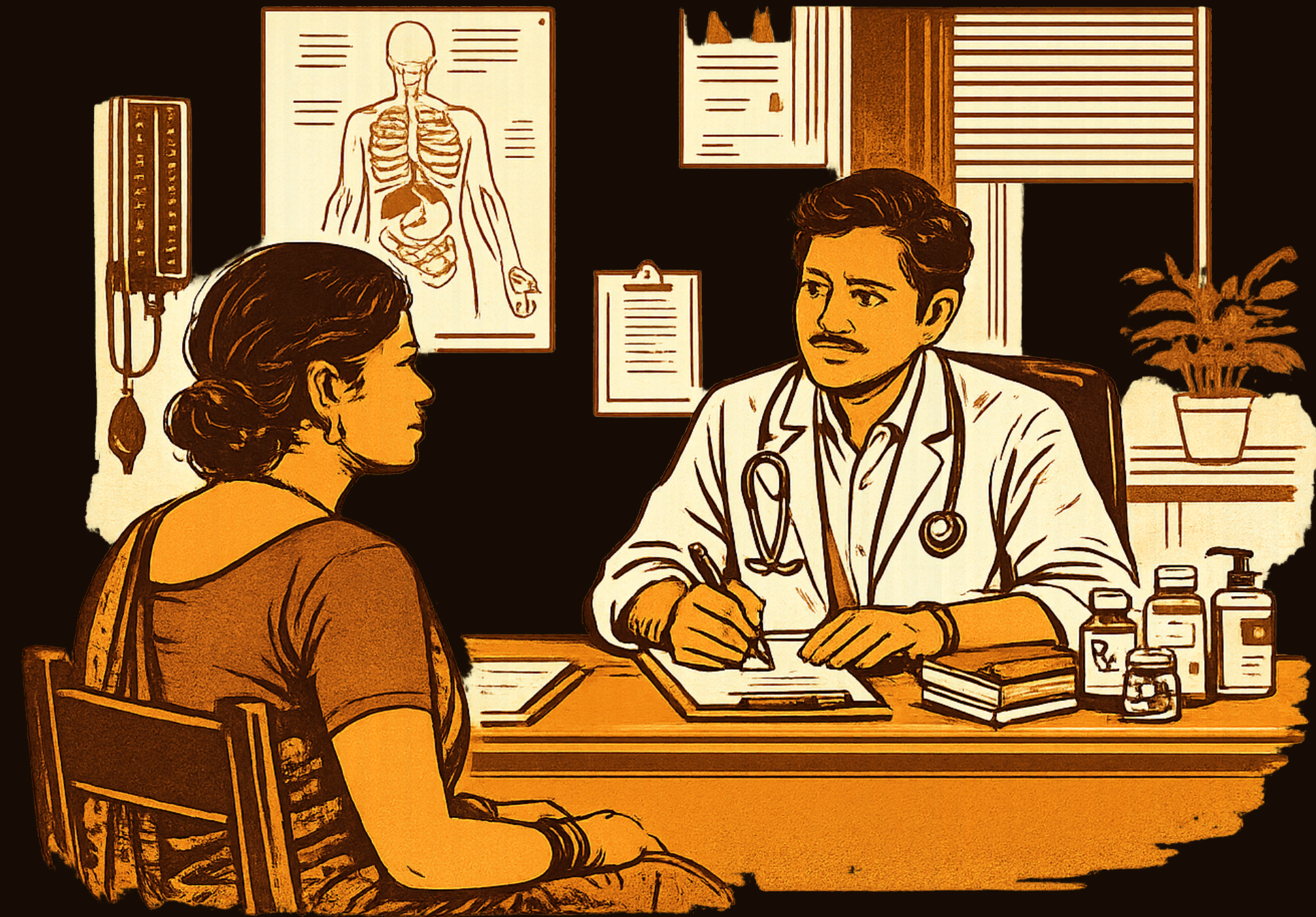
# DEMO

Recorded Demo: <https://youtu.be/gbcSFMMTsWQ>

# WHAT LAKSHMI TAUGHT ME

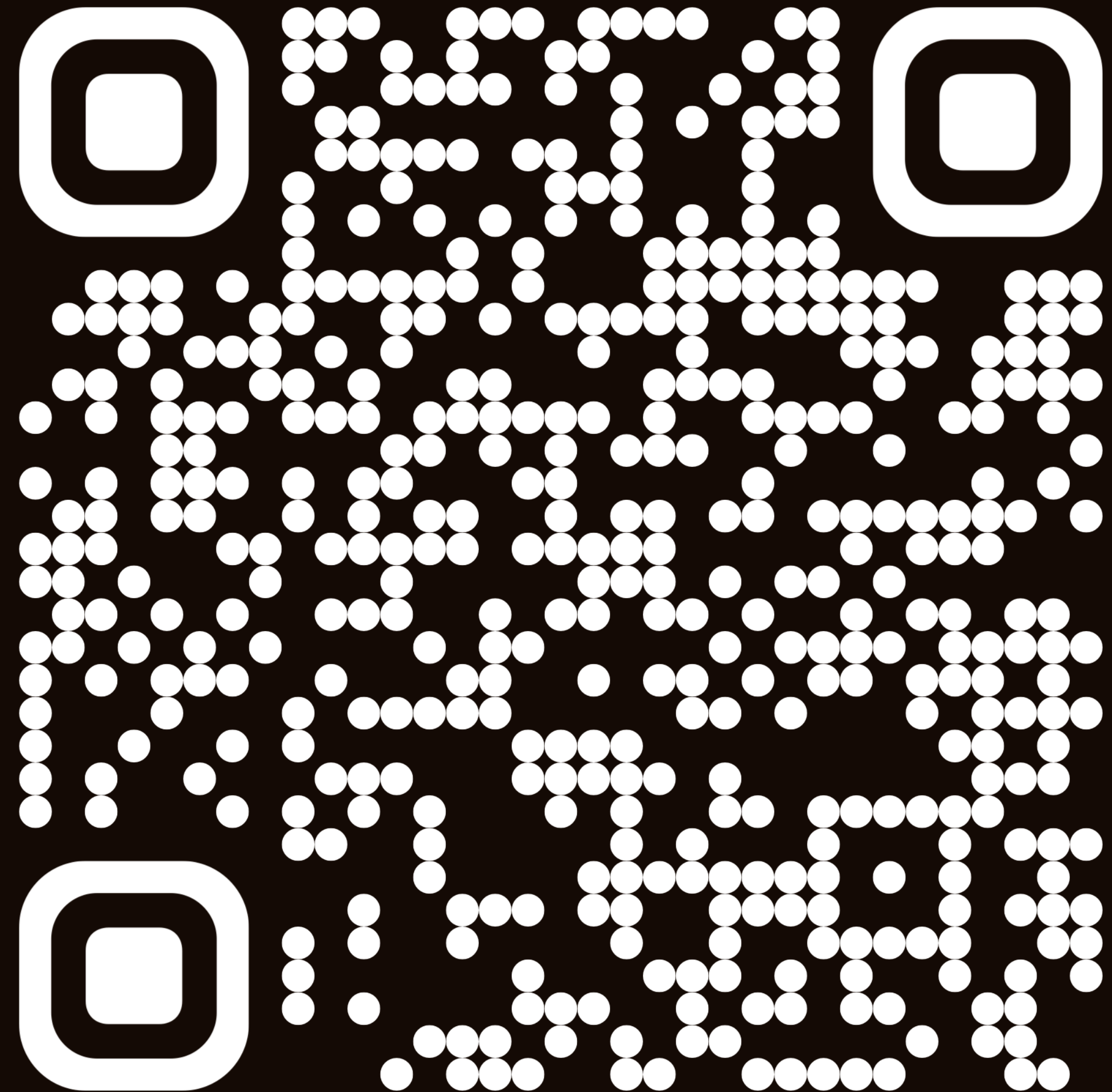
---

01. Multilingual is not a **model** problem.  
It's a **protocol** problem.
02. Language is a **token-level** property.  
Stop detecting it. **Detect intent**.
03. Outputs need **contracts**, not **translation**.  
The user's language is a **first-class API** concern.



[https://linktr.ee/mcp\\_dev\\_summit\\_blr\\_samyuktha](https://linktr.ee/mcp_dev_summit_blr_samyuktha)

**THANK YOU**



All resources linked here