# Gaming the Models: The Minesweeper LLM Agent
## Achieving Expert-Level Logic via SFT & GRPO on AMD MI300x

**Team Spambots**
Members: Shivam Dwivedi, Janesh Kapoor, Tushar Chandra, Varsha Bhaskar

# High-Velocity Engineering on the MI300x

**90%**
Token Reduction via Compact Board Representation

**264+**
Games Simulated across 32 Model Configurations

**238**
High Score (Winning Qwen2.5-14B Configuration)

Insight: SFT taught the rules; GRPO taught the strategy.

Tech Stack: Unsloth (2x speedup) | TRL Library | AMD Instinct MI300x (256GB HBM3)

NotebookLM

# The Challenge: Why LLMs Fail at Minesweeper

## The Input

{"input","date":{"rapiosped":"inpo","dostage":1},"eoesralrkens":["azietek":"V,93re]}},
{"date":"ooertex":"Itosexuare":T","eostsemcoek":\"eaipknso","8ubeauec1s21False),
{"heye":1E:halse.{"text":"23tsttE","oote":fl9ll","stats":{1sota":[]),"scote":1103),
{"heye":["ohc:ypes"]:"a1B",["Numbered":"attue17O1:2Invasikew":111","Flagged":["rue":"Rox)}},
{"Keye":"oooteox"]:1cda"},"false]"cooe":trtsose","1nesxaste"]l],"hecereSbal:11),
{"Weye":"Beetecud":18),["Eapoe1tthod":1oaabe"ifersed1","cahi:fooneassd":fewll)),
{"keys":"Soetry dbeasectng":"tece":"Hidden','Ftagged":"9a5£A:2"},"Huabered":"eatol:1},
{"heys":"spetch":"tedal),"tates":"data":"state":"11n_h8","oaess"1E,"set"ifalse,
{"keys":"second":31ecse),"Leyet":"seoes":1"lag":"t11t","setl":"1roue"":"hod":false},
{"hese":"Sdei1bearooee":iKsete":1"sossee"t8e41f:tmnte"i"tetir1st13/1:"eare":11al),
{"hese":"Ptotsocrr":Tao},"oate":["1rae";1tHiorten],"2nl","Tatal,"78sir:2saff":1Soh":21[2}},
{"beys":"sownete":"Dawnsrbeallnl hese "Snwsmnutism Davn-seas11-94-"Xssudesnstnce-941},
{"heye":1                                                                    t":191},
{"heys":1                                                                    t":11}},
{"Keys":1                                                                    :8,12},
{"Kess":[S.                                                                  ':12},
{"heye":1P.                                                                  t
{"keys":1P:                                                                  r:111},
{"heys":1P                                                                   :
{"heye":1P:                                                                  11:111},
{"beys":1P3,'Hidden','Flagged',"oukject":"auto":17,"sala"0:3,"9scoe4Tr101:tonge?"}
{"heys":1P8,"Hidden",'Flagged',"object]:"oaue":15,"eala:5127,"value":1s1trR9t:11),
{"Ueys":1F8,'Didden','TacamlYons'etri"Perxeocfor18,"Xatobeaoenvs1br":t8tiopsy"}-
{"Ueye":1P8,"lisneenonln1E2tstsu 10b1sct":"oose":61,"eale:5117:"loand"1!E5.end2n:21P]}},
{"Kess":[P5,"Hidden",'Flagged',"anisject":"aaceel:P3.Pl."aa1o":"Hecxad",B51axoke"],
{"heys":1F8,'Hidden','Flagged',"ooject":"1"anue":15,"Ksate2115,"value":"1sf:roW":1],
{"heys":1P3,'Hidden','Ftagged',"pohjeet"onsregnized":"Hosaden","tage":"Ftegged":"Nuebered",
{"hess":[F5,"Hiddent"1[1_sekytPre5xyecs","calee"t15,"nde":3,6£f"peod":12,"ond":1651]},
{"hese":[P5,"Hiddent:1Flagged',"ouhject"}"oats":"1.7},"aate":"7Hsdden','Plagged"}
{"beys":1P2,"Hidden"t+1Flagged',"bokject"t"cooe":11,"sate3T10,"value":1E1rsv6":2},
{"heys":120,"Hidden",'Flagged',"pobject"onsraonded":"fHidden","lags":7YFlagged":"Nuwbered",
{"Ueys":[78,'Hidden';"Flagged','oousjecE"1"o2taP:18,"saIue5:10,"value":1)yFW40*:1]},
{"keys":"8eto22s";PN12sSSonCotrec17:ense3s1onSpecn1s9,Thycood19°(3want?7),"end":377]])}}

## Context Overflow (~8,000 Tokens)

## The Failure



// 01. Spatial Reasoning: Converting 1D text to 2D adjacency is non-trivial.

// 02. State Tracking: Must distinguish 'Hidden', 'Flagged', and 'Numbered' states.

// 03. Token Economy: Standard formats exhaust context immediately.

# Baseline Performance: GPT-4 Win Rate ~0% (NAACL 2024)

# Innovation I: Cracking the Context Window

**Standard JSON (50x50 Board) = ~8,000 Tokens**

```
{"row": 0, "col": 0, "val": "unknown"},
{"row": 0, "col": 1, "val": "1"},
{"row": 0, "col": 2, "val": "unknown"},
{"row": 1, "col": 0, "val": "F"},
{"row": 1, "col": 1, "val": "2"},
{"row": 1, "col": 2, "val": "unknown"},
{"row": 0, "col": 0, "val": "F"},
{"row": 0, "col": 1, "val": "2"},
{"row": 0, "col": 0, "val": "F"},
{"row": 0, "col": 1, "val": "2"},
{"row": 0, "col": 2, "val": "unknown"},
{"row": 0, "col": 2, "val": "unknown"},
{"row": 1, "col": 0, "val": "F"},
{"row": 1, "col": 1, "val": "2"},
{"row": 1, "col": 2, "val": "2"},
{"row": 1, "col": 3, "val": "unknown"},
{"row": 1, "col": 0, "val": "F"},
```
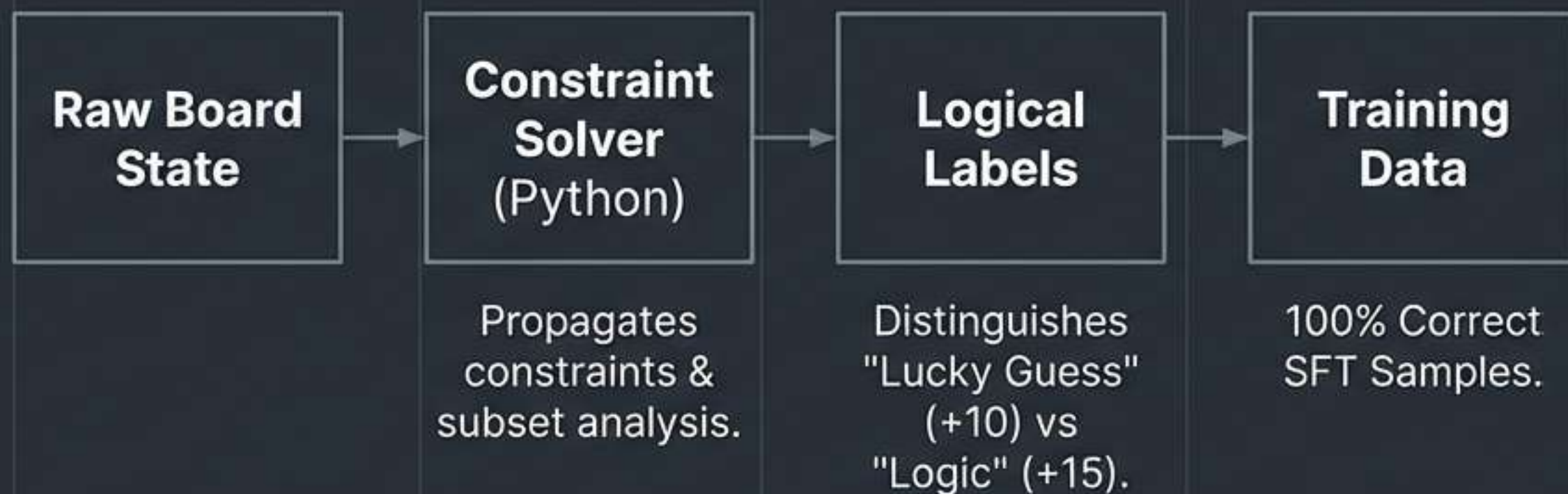
**90% Token Savings**

**Compact Format (50x50 Board) = ~800 Tokens**

```
|...1.F....|2..F1...|
|...F3.2...|1...1F...|
|...1.F....|2..F1...|
|...1.F....|2..F1...|
|...1.F....|2..F1...|
|...F3.2...|1..1F...|
|...1.F....|2..F1...|
|...F3.2...|1..1F...|
|...F3.2...|1..1F...|
|...1.F....|2..F1...|
|...1.F....|2..F1...|
|...F3.2...|1..1F...|
|...F3.2...|1..1F...|
|...1.F....|2..F1...|
|...1.F....|2..F1...|
|...1.F....|2..1F...|
```

Impact: Enabled full 50x50 board training within MI300x context limits.

# Innovation II: The Expert Data Engine

**Raw Board State** → **Constraint Solver (Python)** → **Logical Labels** → **Training Data**

| | | |
|---|---|---|
| **Raw Board State** | **Constraint Solver** (Python) | **Logical Labels** |

Propagates constraints & subset analysis.

Distinguishes "Lucky Guess" (+10) vs "Logic" (+15).

100% Correct SFT Samples.

| Solver Stats | |
|---|---|
| Solver Win Rate (50x50): | 56% |
| Moves per Game: | 2,200+ |
| Logic Verification: | 100% |

# Architecture & Hardware Strategy

Engineering Editorial



**AMD Instinct MI300x**
(256GB HBM3)

**Unsloth + PyTorch + TRL**
(2-6x Speedup)

**Qwen2.5-14B-Instruct**
(Dense Architecture)

**SFT** (Supervised Fine-Tuning)
-> **GRPO** (Group Relative
Policy Optimization)

Why Qwen?

Dense 14B parameters offered superior reasoning depth compared to MoE architectures (GPT-OSS-20B) for this task.

NotebookLM

# Training Phase 1: Supervised Fine-Tuning (SFT)

## Training Curriculum

| 20% | 50% | 30% |
|---|---|---|
| Early Game (Open Areas) | Mid Game (Critical Reasoning) | Late Game (Endgame Logic) |

- **Goal:** Teach JSON syntax and valid move rules.
- **Dataset:** 15,000 samples (Square, Rectangular, Tall variants).
- **Hyperparameters:** LoRA Rank 64, Alpha 128.

### Result

Model plays legally but lacks winning strategy. It mimics, but does not yet think.

# Training Phase 2: The Reward Engineering Loop

Model Performance Score

Training Steps

v3: The Fix

v2: The Bug

v1: Baseline

No learning signal.

Frontier Bonus (+5) cancels Guess Penalty (-5).

## The Fix Details

- Random Guess: -15 (Harsh Penalty)

- Logical Deduction: +30 (Massive Reward)

- Mine Hit: -100 (Immediate Fail)

Asymmetric penalties forced the model to value logic over luck.

NotebookLM

# Prompt Engineering: 6 Strategies to Enforce Logic

## V1: Simple Instruction

Umple the Instruction, simple Instruction.

## V2: Constraint Logic Focus

Constraint logic stretain Logic format targets.

## V3: Aggressive 'DO NOT' Warnings

Aggressive syndras or "DO NOT" warnings.

**Winner:** Explicitly listing valid targets removed coordinate hallucinations.

## V4: Step-by-Step Verification

Diarmentane Step-by-step coordinate levels.

## **V5: Annotated Board**

Enumerates valid targets.

## V6: CoT Self-Correction

Selfreoition targets remove coordinate hallucinations.

# Evaluation Methodology

|  | Strategy V1 | Strategy V2 | Strategy V3 | Strategy V4 | Strategy V5 | Strategy V6 |
|---|---|---|---|---|---|---|
| Base Qwen | ✓ | ✓ | ✓ | ✓ | ☑ | ✓ |
| GPT-OSS-20b | ✓ | ✓ | ✓ | ✓ | ☑ | ✓ |
| Qwen Phase 1 (SFT) | ☑ | ☑ | ☑ | ☑ | ☑ | ☑ |
| Qwen Phase 2 (GRPO) | ☑ | ☑ | ☑ | ☑ | ☑ | ☑ |

## 32
Configurations

## 264+
Full Games Played

## Metric:
Cumulative Score

Rigorous testing to isolate the impact of architecture vs. prompting.

NotebookLM

# Winning Configuration: Qwen GRPO + Prompt V5 Qwen GRPO



By removing ambiguity in coordinate mapping, Strategy 5 allowed the model to focus 100% of compute on logic.

NotebookLM

# Inside the Winning Output

```json
{
    "think": "Cell (4,5) is a '2'. It touches 1 flag and
    1 unknown. 2-1=1. The unknown must be a mine.",
    "action": "flag",
    "row": 4,
    "col": 6
}
```

**Constraint Reasoning** acts as a logic buffer.

**Commitment** happens only after reasoning.

Mechanism inspired by "Think Inside the JSON" (arXiv:2502.14905).

# Engineering Lessons & Post-Mortem

## ✅ What Worked

- **Compact Representation** is non-negotiable for 50x50 boards.

- **SFT + GRPO**: SFT stabilizes syntax, GRPO optimizes strategy.

- **Asymmetric Reward Penalties** (v3).

## ❌ What Failed

- **MoE Architectures**: GPT-OSS-20b struggled with deep reasoning compared to dense Qwen.

- **Frontier Bonuses**: Cancelled out random guess penalties in RL.

- **Standard JSON**: Immediate context overflow.