

Lessons Learned: A Multi-Agent Framework for Code LLMs to Learn and Improve

Yuanzhe Liu¹, Ryan Deng², Tim Kaler², Xuhao Chen^{2,3},
Charles E. Leiserson², Yao Ma¹, Jie Chen⁴

¹Rensselaer Polytechnic Institute ²Massachusetts Institute of Technology

³Michigan State University ⁴MIT-IBM Watson AI Lab, IBM Research


NeurIPS 2025

Motivation

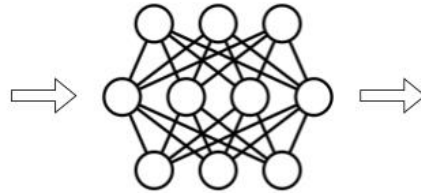
- Code Optimization is less explored.
- Task of code optimization: $f(\text{slower-code}) = \text{faster-code}$
- Slower code is
 - Compilable & Correct

```
#include <iostream>
using namespace std;

int main(){
    int n;
    cin >> n;
    int sum = 0;
    for (int i = 1; i <= n; i++) {
        sum += i;
    }
    cout << sum << endl;
    return 0;
}
```




(a) Slower Code.



```
#include <iostream>
using namespace std;

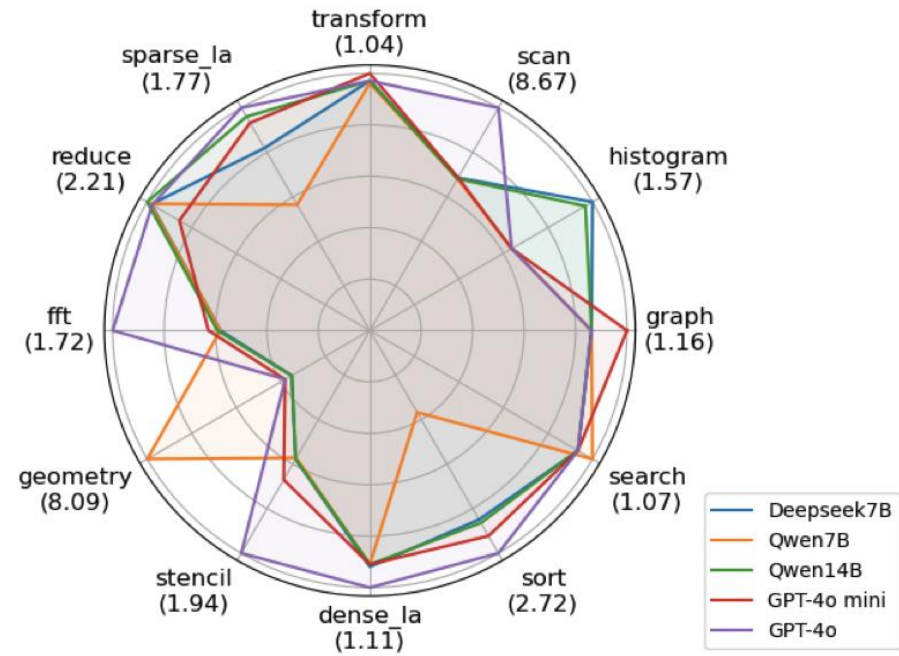
int main(){
    int n;
    cin >> n;
    cout << n*(n+1)/2 << endl;
    return 0;
}
```



(b) Faster Code.

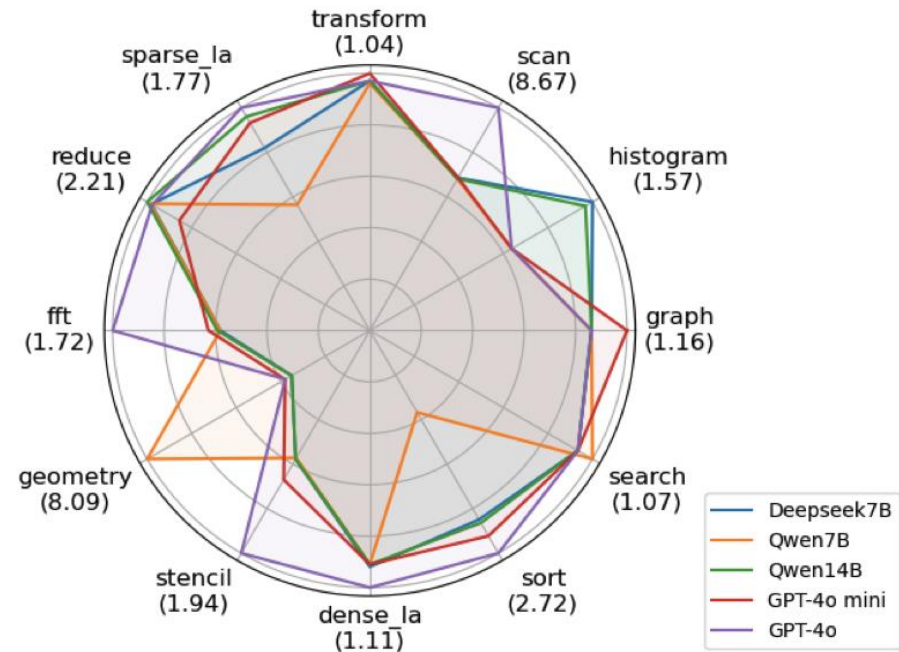
Our findings

- On ParEval benchmark, no one LLM performs the best on all problems



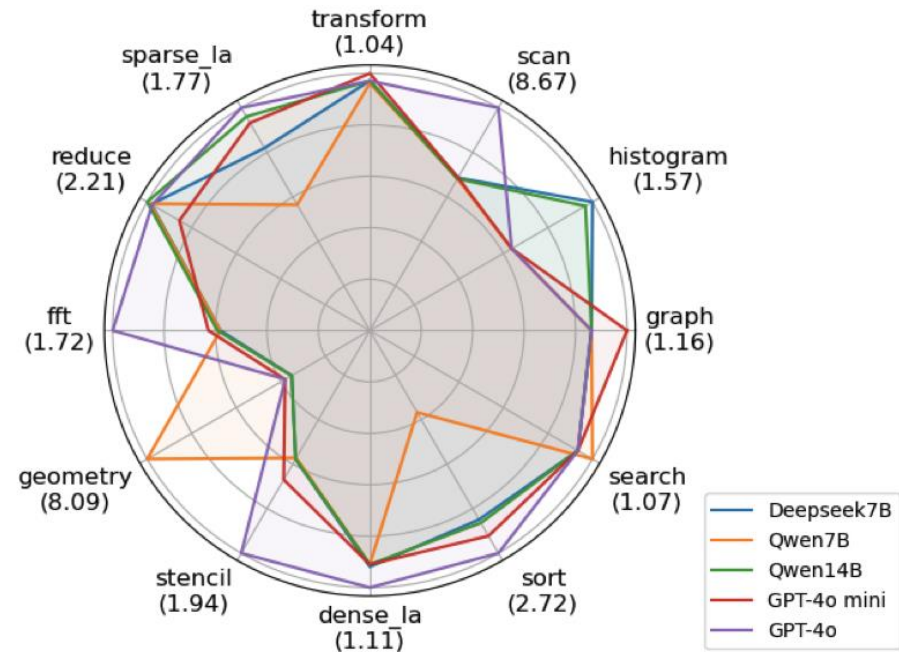
Our findings

- On ParEval benchmark, no one LLM performs the best on all problems
- GPT-4o is the overall winner.



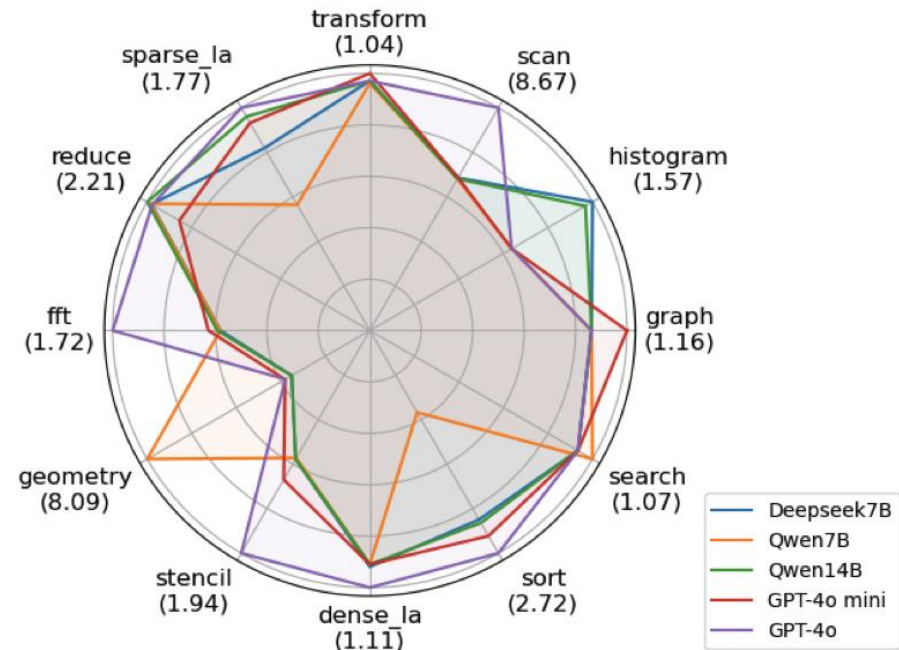
Our findings

- On ParEval benchmark, no one LLM performs the best on all problems
- GPT-4o is the overall winner.
- E.g. on “Histogram” Deepseek7B and Qwen14B outperform GPT-4o by $1.6\times$



Our findings

- On ParEval benchmark, no one LLM performs the best on all problems
- GPT-4o is the overall winner.
- E.g. on “Histogram” Deepseek7B and Qwen14B outperform GPT-4o by $1.6\times$
- Different LLM presents distinct capabilities.



Usage of Multiple LLMs

- How to use multiple agents to solve a coding problem?

Usage of Multiple LLMs

- How to use multiple agents to solve a coding problem?
 - We advocate the concept of *lessons*

Usage of Multiple LLMs

- How to use multiple agents to solve a coding problem?
 - We advocate the concept of *lessons*
 - Such *lessons* are summarized by LLM agents, learned by others, so they can collectively improve the code performance.

Usage of Multiple LLMs

- How to use multiple agents to solve a coding problem?
 - We advocate the concept of *lessons*
 - Such *lessons* are summarized by LLM agents, learned by others, so they can collectively improve the code performance.

Original code

```
for (int i = 0; i < n; ++i)
  for (int j = 0; j < n; ++j)
    for (int k = 0; k < n; ++k)
      C[i][j] += A[i][k] * B[k][j];
```

Naive implementation of matrix multiplication $C = AB$.

Improved code, round 1

```
for (int i = 0; i < n; ++i)
  for (int k = 0; k < n; ++k)
    for (int j = 0; j < n; ++j)
      C[i][j] += A[i][k] * B[k][j];
```

Lesson: Reordering loops improves cache locality and increases performance. The order of (i, k, j) out of 6 different permutations often performs the best, because of how caches work.

Usage of Multiple LLMs: LessonL

How do agents work together in practice by *lesson*?

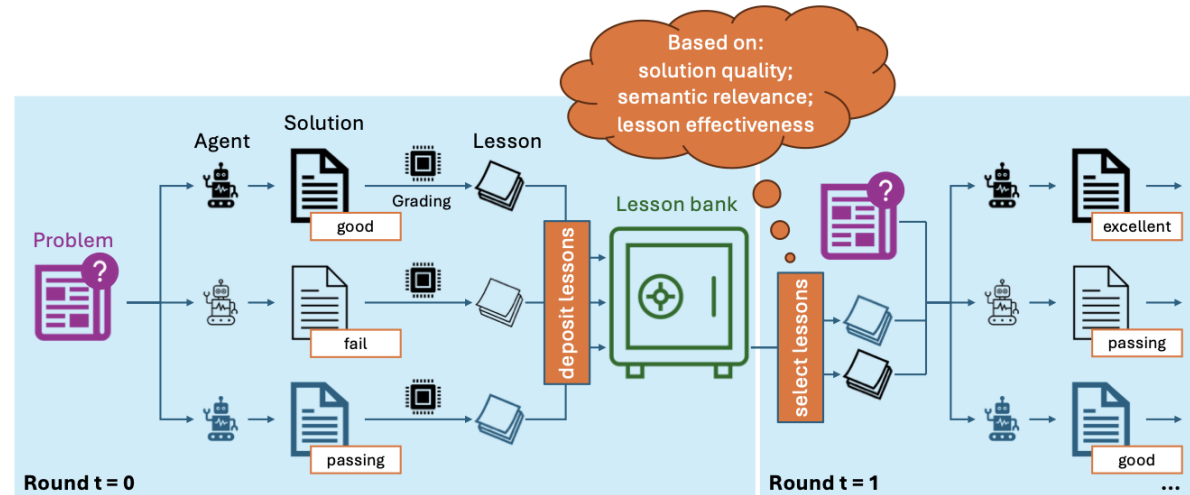


Figure 2: The LessonL framework (which may repeat multiple rounds).

Usage of Multiple LLMs: LessonL

How do agents work together in practice by *lesson*?

- Lesson Solicitation
 - under different scenarios

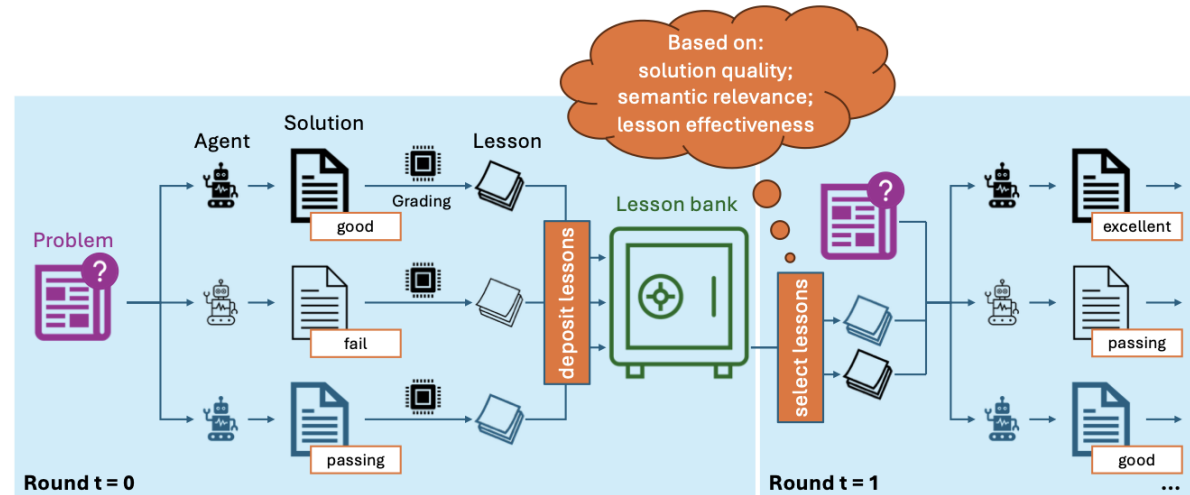


Figure 2: The LessonL framework (which may repeat multiple rounds).

Usage of Multiple LLMs: LessonL

How do agents work together in practice by *lesson*?

- Lesson Solicitation
 - under different scenarios
- Lesson Banking and Selection
 - reduce context length

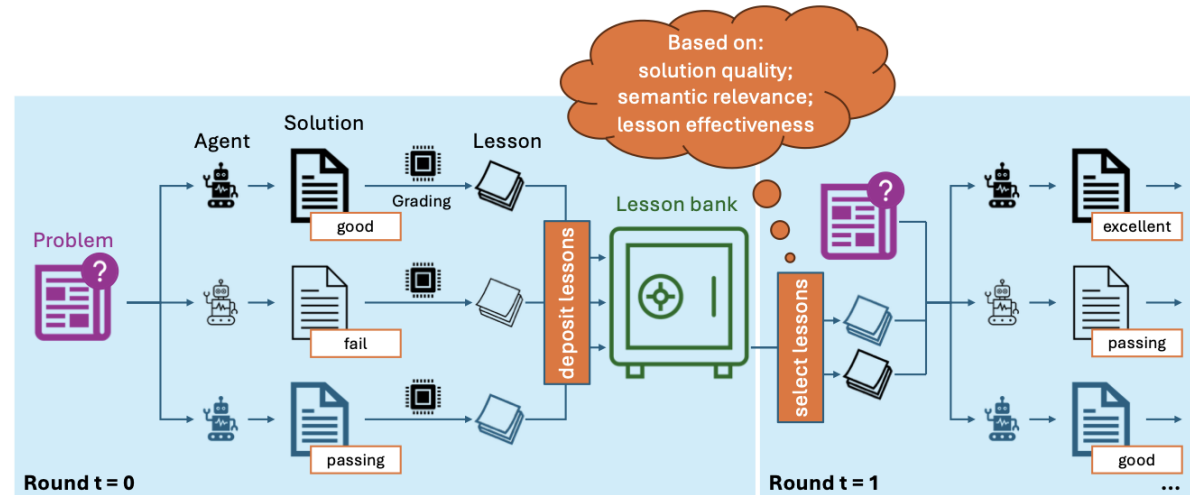


Figure 2: The LessonL framework (which may repeat multiple rounds).

Usage of Multiple LLMs: LessonL

How do agents work together in practice by *lesson*?

- Lesson Solicitation
 - under different scenarios
- Lesson Banking and Selection
 - reduce context length
- Effectiveness Adjustment factor f :
 - adjust selection dynamically

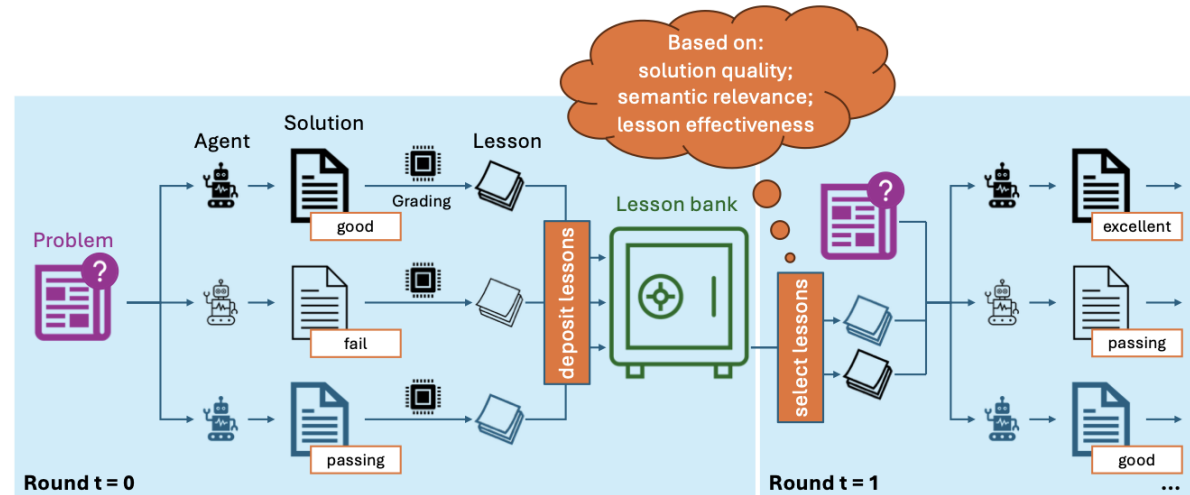


Figure 2: The LessonL framework (which may repeat multiple rounds).

Usage of Multiple LLMs: LessonL

How do agents work together in practice by *lesson*?

- Lesson Solicitation
 - under different scenarios
- Lesson Banking and Selection
 - reduce context length
- Effectiveness Adjustment factor f :
 - adjust selection dynamically

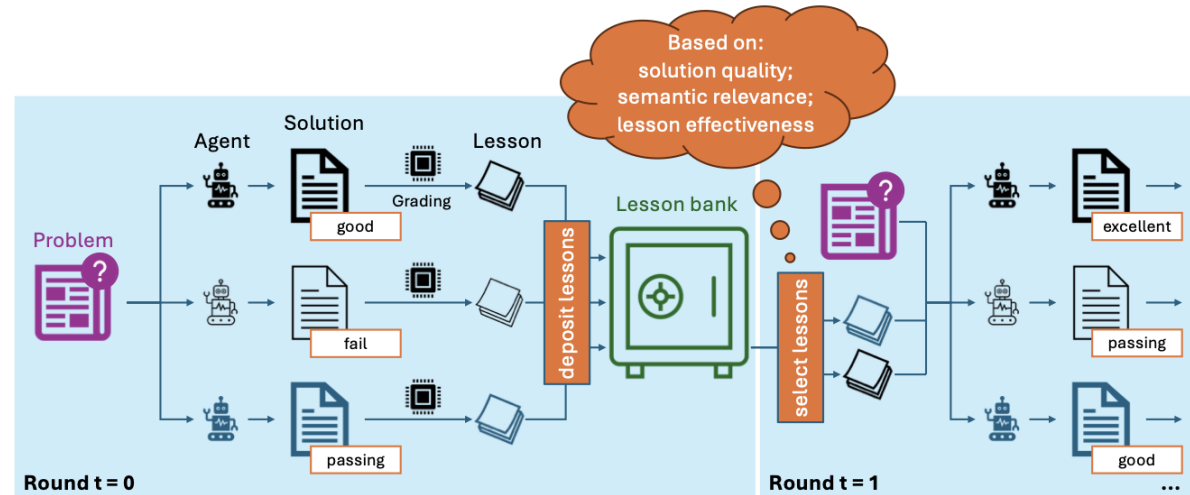


Figure 2: The LessonL framework (which may repeat multiple rounds).

Results on Code Optimization

ParEval	Serial mode			OpenMP mode		
	Correct	> 2x	Speedup	Correct	> 2x	Speedup
GPT-4o	0.80 ± 0.00	0.16 ± 0.03	1.72 ± 0.11	0.73 ± 0.05	0.58 ± 0.05	2.93 ± 0.30
OpenAI o3	0.77 ± 0.02	0.23 ± 0.04	2.21 ± 0.16	0.72 ± 0.03	0.58 ± 0.03	3.55 ± 0.27
MapCoder	<u>0.88 ± 0.02</u>	0.15 ± 0.02	1.85 ± 0.08	<u>0.83 ± 0.05</u>	0.58 ± 0.02	3.43 ± 0.17
LessonL	<u>0.91 ± 0.02</u>	<u>0.21 ± 0.01</u>	<u>2.16 ± 0.11</u>	<u>0.86 ± 0.01</u>	<u>0.62 ± 0.02</u>	<u>3.46 ± 0.03</u>

LessonL models:

- deepseek-coder-7b-instruct-v1.5
- Qwen2.5-Coder-7B-Instruct
- Qwen2.5-Coder-14B-Instruct

Thanks for listening!

NeurIPS 2025 Exhibition Hall C,D,E
Thu 4 Dec 4:30 p.m. – 7:30 p.m. PDT



paper



code