# TOWARDS AI-NATIVE SOFTWARE ENGINEERING (SE 3.0): A VISION AND A CHALLENGE ROADMAP

Presenter: Haonan Zhang January 16, 2025



# Different stages of SE



activities (AI4SE)

• Powered by expensive

NOW

data-driven models with limited capabilities



development

 Tools supporting traditional SE Process activities

• Powered by program

analysis technologies

PAST

ig • Al-native SE process ss maximizing the strengths of human (requirements) & Al (implementation)

Powered by efficient
knowledge-driven

knowledge-driven models with advanced reasoning capabilities

time

ming (ICML'24)

4

2

reasoning capabi

#### Significant cognitive overload on the human



3

#### Inefficient model training on domain-specific tasks

An example of the performance of FMs used for domain-specific tasks

C	Cross-Task		Cross-Website			Cross-Domain		
Ele. Acc	Op. F1	Step SR	Ele. Acc	Op. F1	Step SR	Ele. Acc	Op. F1	Step SR
ne-Tuning								
40.5	74.4	37.3	28.7	69.6	27.9	38.2	69.1	36.2
52.2	70.7	48.8	35.3	65.8	32.7	41.9	64.6	39.5
56.8	74.6	52.5	42.6	69.9	39.5	43.8	65.2	40.7
39.5	74.9	36.1	34.0	70.8	32.2	38.2	72.8	37.5
50.0	72.1	46.0	39.5	71.5	36.3	40.9	70.1	39.4
52.9	74.9	50.3	41.7	74.1	38.3	43.8	73.4	39.6
arning								
19.4	59.8	16.8	14.9	56.5	14.1	25.5	57.9	24.2
40.2	63.4	31.7	27.4	61.0	27.0	36.2	61.9	29.7
4.7	39.5	4.7	9.7	37.8	9.7	16.0	41.4	15.3
15.1	66.5	13.0	11.3	63.4	10.5	16.5	65.1	14.7
48.9	69.1	40.6	48.5	70.6	41.7	44.0	70.9	40.9
72.9	80.9	65.7	74.4	83.7	70.0	72.8	73.6	62.1
	40.5 55.2 56.8 39.5 50.0 52.9 arning 19.4 40.2 4.7 15.1 48.9 72.9	Cross-Tas       Ele. Acc     Op. F1       e-Tuning     40.5       40.5     74.4       52.2     70.7       56.8     74.6       50.0     72.1       52.9     74.9       arring     19.4       40.2     63.4       40.2     63.4       4.7     39.5       15.1     66.5       48.9     69.1       72.9     80.9	Cross-Task       Ele. Acc     Op. Fl     Skep SR       ser-Tuning     40.5     74.4     37.3       52.2     70.7     48.8     65.8       30.5     74.9     36.1     50.2       50.2     74.9     50.3     30.1       10.4     59.8     16.8     16.8       40.2     63.4     31.7     17.1       47.7     39.5     4.7     30.5       15.1     66.5     13.0     48.9     69.1       47.9     69.3     40.0     65.7     14.9	Cross-Task     Cro       Ele. Acc Op. FI Step SR     Ele. Acc       ex-Tuning     105     74.4     37.3     28.7       62.0     70.7     48.8     35.3     56.8     74.6     35.3       50.8     74.0     36.1     34.0     30.5     72.1     46.0     39.5       25.9     74.9     36.1     34.0     30.5     52.4     24.6       30.5     72.1     46.0     39.5     50.3     1.17     rarring       10.4     59.8     16.8     14.0     31.7     27.4       4.7     39.5     4.7     9.7     13.0     11.3       11.7     15.1     66.5     13.0     11.3     48.9     05.7     74.4       7.2     9.0     9.1     40.0     48.5     77.4     9.5     74.9	Cons-Task     Cross-Web       Ele. Acc 0p, Fl Step SR     Ele. Acc 0p, Fl       Mark 200, Fl Step SR     Ele. Acc 0p, Fl       Mark 201, Fl Arrow 201, Fl Arro	Cross-Task     Cross-Website       Ele. Acc     Op. F1     Step 5R     Ele. Acc     Op. F1     Step 5R       en-Tuning     40.5     74.4     37.3     28.7     69.6     27.9       52.2     70.7     48.8     35.3     65.8     32.7     59.6     27.9     36.1     34.0     70.8     32.2     50.0     71.1     83.3     36.5     37.9     36.3     37.5     36.3     37.7     36.4     36.5     14.1     40.2     56.4     31.6     37.7     37.4     36.3     37.0     37.6     37.6     37.6     36.4	Cross-Task     Cross-Website     Cross-Website       Tel. Acc Op, FI     Step SR     Ele. Acc     Op, FI     Step SR     Ele. Acc       en-Tuning     40.5     74.4     37.3     28.7     60.6     27.9     38.2       92.2     70.7     48.8     33.3     65.8     32.7     41.9       92.9     74.9     36.1     34.0     70.8     32.2     38.2       92.9     74.9     36.1     34.0     71.5     36.3     40.3       92.9     74.9     36.1     34.0     71.4     38.3     43.8       92.9     74.9     36.1     34.0     70.8     32.2     38.2       92.9     74.9     36.1     34.0     70.5     36.3     40.3       92.9     74.9     36.1     34.0     71.4     38.3     43.8       string     50.3     41.7     74.1     38.3     43.8       string     50.3     41.7     74.1     38.3     43.8       10/4     <	Cons-Task     Cross-Website     Cross-Dom       Ele. Acc     0p, Fl     Sep SR     Ele. Acc     0p, Fl     Step SR     Ele. Acc     0p, Fl

"The most effective grounding strategies we explored in this paper still exhibit a 20-25% performance gap compared to oracle grounding"

Boyuan Zheng, Boyu Gou, Jihyung Kil, Huan Sun, and Yu Su. 2025, GPT-4V(ision) is a generalist web agent, if grounded. It

#### Unaffordable resource usage

It is very expensive to fine-tune the models for improved domain-specific capability

"We predict that to achieve 95% accuracy for in-domain low-level tasks, 1M episodes would be required, while 2M episodes would be required to obtain 95% episode completion rates for 5-step high-level tasks"

Li. Wei, et al. On the Effects of Data Scale on Computer Control Agents. In Proceedings of the 38th Conference on Neural Information Processing Systems (2)

5

7

#### Suboptimal code quality produced by FMs



"While functionally correct, this approach suffers from sub-optimal time complexity and space complexity" Nog HUNK, Value (NN, Weij Shang, Hening Cal, and Je Zang, Efflench: Rendmarking the Efficiency of Automatically Conversed Code In Proceedings of the 5th Conference on Neural Information Proceeding System (NeurIPS2)

6

8

#### **Vision of Software Engineering 3.0**



Development is no longer driven by code but by intents expressed through backand-forth conversations between human developers and their AI teammates.

#### What fuzzing looks like in SE 3.0





#### What test execution looks like in SE 3.0



٩

11

#### What GUI testing looks like in SE 3.0



Q1 (Start prompt): We want has following activities... The operation is required? (<Op uppe Wha

What GUI testing looks like in SE 3.0

#### Video Demo

Source: https://osu-nlp-group.github.io/SeeAct/

Balance between ask too many and not asking enough

**Optimizing Instructions and Demonstrations** for Multi-Stage Language Model Programs

Krista Opsahl-Ong<sup>1</sup>\*, Michael J Ryan<sup>1</sup>\*, Josh Purtell<sup>2</sup>, David Broman<sup>3</sup>, Christopher Potts<sup>1</sup>, Matei Zaharia<sup>4</sup>, Omar Khattab<sup>1</sup>

<sup>1</sup>Stanford University, <sup>2</sup>Basis, <sup>3</sup>KTH Royal Institute of Technology <sup>4</sup>UC Berkeley



Figure 1: An example of the optimization problem we explore, shown for a multi-hop retrieval LM program. Given some question-answer pairs and a metric, the optimizer proposes new instructions and bootstraps new demonstrations (not pictured) for each stage.

10

### Improving the efficiency of code synthesis

#### SynCode: LLM Generation with Grammar Augn Shubham Ugare Daternity of Elizate Urbane Champaign, 553. Tarun Suresh tena Champaign, 55A Hangoo Kong University of Elizate Urbano Champaijas, 8334 Sasa Misalkwie Gagandeep Singh University of Elizate I I M Completed Code C

Figure 1: In the SynCode workflow, the LLM takes partial output  $C_k$  and generates a distribution for the next token  $t_{k+1}$ . The parser processes  $C_k$  to produce accept sequences A and remainder r. These values are used by the DFA mask store to create a token mask, eliminating syntactically invalid tokens. The LLM iteratively generates a token  $t_{k+1}$  using the distribution and the mask, appending it to  $C_k$  to create the updated code  $C_{k+1}$ . The process continues until the LLM returns the final code  $C_n$  based on the defined stop condition.

#### Improving runtime performance



**Negative points** 

13

15

- > Mostly from a theoretical perspective...
- ➤ Blurry boundary between SE 2.0 and SE 3.0?

the human develope". In our experience, a typical programming session looks as follows: create a class, write the constructor's signature and have the copilot autocomplete the implementation, create a new method, write a code comment inside the body of that methods of  $(\pi_{12}, \pi_{22}, \pi_{22}$ 



14

16

# **Positive points**

- > Redefine software engineering in the era of LLM
- > Reveal the challenges we need to address to shift to SE 3.0.

Source: https://osu-nlp-group.github.io/SeeAct

#### Rating

Theoretical foundation: 5/5

Practical experience: 4/5

Overall: 4.5/5

#### Discussion

> Is SE 3.0 = SE based on agent?

#### 'Modern' agent = LLM + external environment?



Source: https://ysu1989.github.io/resources/language\_agents\_YuSu\_2024.pdf

#### Discussion

> We are now at SE 2.0 or 2.5?



#### Discussion

17

19

> Do FMs really think or just content retrieval?



Source: https://x.com/

> If just content retrieval, then can we ever shift to SE 3.0?

18

## Discussion

> Is Prompt Engineering the only thing we can do as a SE researcher?





22

Source: https://x.com/ 21