# Software Engineering

# (9) State-of-the-Art Topics

Sokendai / National Institute of Informatics

Fuyuki Ishikawa / 石川　冬樹

f-ishikawa@nii.ac.jp / @fyufyu

http://research.nii.ac.jp/~f-ishikawa/

大学共同利用機関法人 情報・システム研究機構
国立情報学研究所
National Institute of Informatics

# TOC

- **<u>AI for SE</u>**
- SE for AI
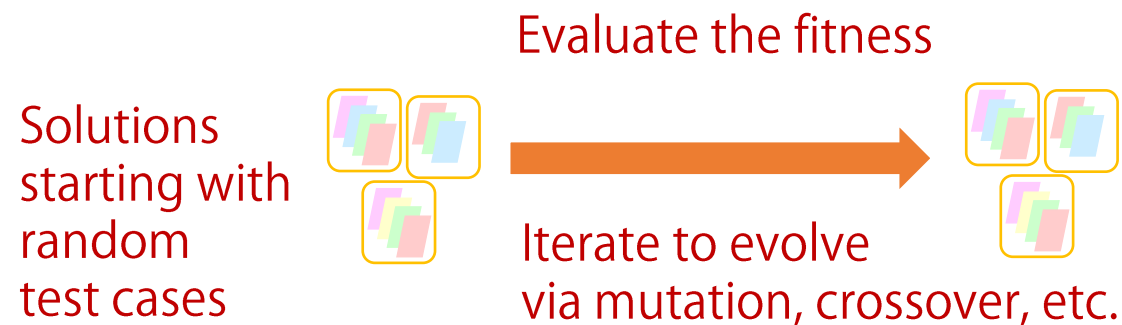- Generative Chat AI

# AI for SE?

- New directions appeared with the AI trend with deep learning in late 2010's
  - e.g., bug prediction of program code
- Advanced automation techniques had been investigated a lot before it
  - e.g., use of metaheuristics optimization or evolutionary computation such as genetic algorithms has been popular for these 10 years (called "search-based software engineering")
  - Showing some of them (maybe biased)

# Active Research Area (1) Search-Based SE

- Search-Based Software Engineering:
  - Solve problems in software engineering by reducing them to optimization problems, i.e., define some "score" and maximize it
  - Often done with metaheuristic algorithms
  - Test generation, test minimization, test prioritization, design, requirements selection, bug fix creation, …
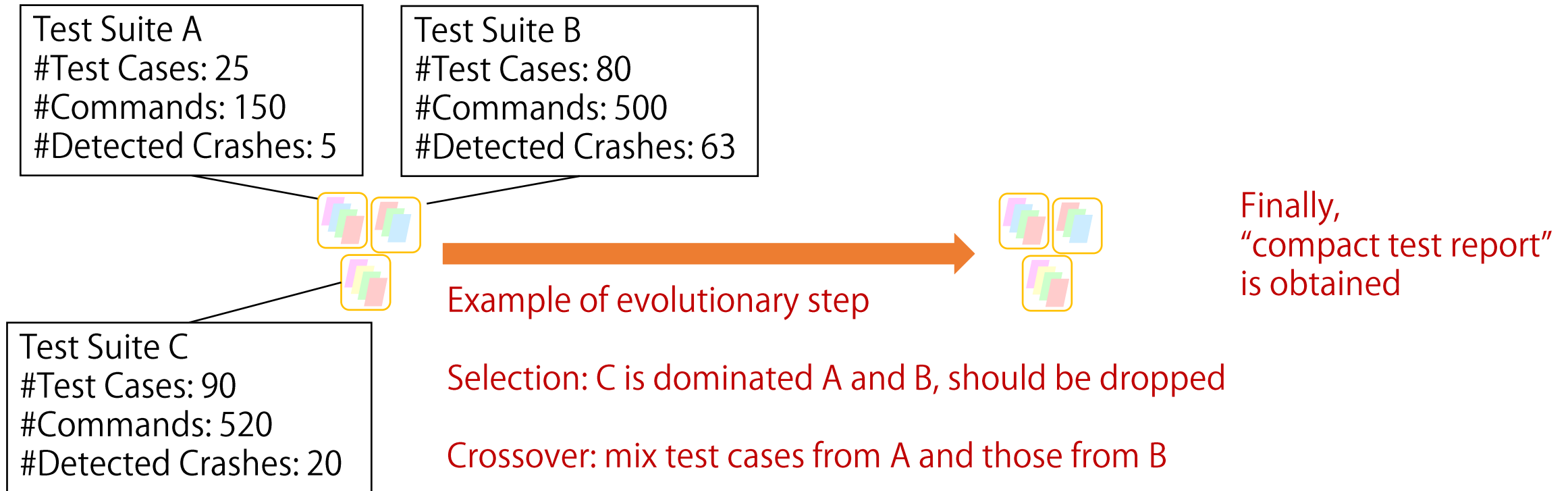
# Example: Search-Based Testing

- Example: automated test input generation
  - For the unit testing of the method `m1(int x, int y)`, prepare the input vector `(x0, y0), (x1, y2), …, (xn, yn)`
  - Maximize the branch coverage, while minimizing the number of test cases `n`

Evaluate the fitness

Solutions starting with random test cases

Iterate to evolve via mutation, crossover, etc.

[ https://www.evosuite.org/ ]
[ Fraser et al., Whole Test Suite Generation, TSE'13 ]

# Application Example: Mobile App Testing at Facebook

■ Sapienz: automated test generation

Test Suite A
#Test Cases: 25
#Commands: 150
#Detected Crashes: 5

Test Suite B
#Test Cases: 80
#Commands: 500
#Detected Crashes: 63

Test Suite C
#Test Cases: 90
#Commands: 520
#Detected Crashes: 20

Finally,
"compact test report"
is obtained

Example of evolutionary step

Selection: C is dominated A and B, should be dropped

Crossover: mix test cases from A and those from B

[ Mao et al., Sapienz：multi-objective automated testing for Android applications, ISSTA'16 ]
[ Alshahwan., et al., Deploying Search Based Software Engineering with Sapienz at Facebook, SSBSE'18 ]

# Active Research Area (2) Fault Localization

■Fault Localization: estimating the bug location(s)
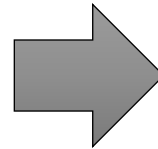
■Example: Spectrum-based Fault Localization

→ Compare lines executed by passed tests and failed tests

Test record

| l.1 | l.2 | l.3 | ⋯ | Result |
|-----|-----|-----|-----|--------|
| ✓ | ✓ | | ⋯ | PASS |
| ✓ | | ✓ | ⋯ | FAIL |
| | | ✓ | ⋯ | FAIL |
| ⋯ | ⋯ | ⋯ | ⋯ | ⋯ |

"Suspiciousness" ranking

| Line | Score |
|------|-------|
| l.3 | 0.8 |
| l.5 | 0.72 |
| l.9 | 0.6 |
| ⋯ | ⋯ |

$$\text{Ochiai Metric (line)} = \frac{\#\ Failed\ Test\ Cases\ that\ execute\ the\ line}{\sqrt{\#\ Test\ Cases\ that\ execute\ the\ line \times \#\ Failed\ Test\ Cases}}$$

[ Wong et al., A Survey on Software Fault Localization, TSE'16 ]

# Active Research Area (3) Automated Program Repair

■ <span style="color:red">Automated Program Repair</span>

  ■ Given a test suite and a program, modify the program so that it passes all of the test cases in the test suite

  ■ Common approach: "generate and validate"

    ■ e.g., apply (inverse) mutation operators that denote typical bugs, such as changing between < and <=

    ■ e.g., search for the best operators by evolutionary computation

  ■ "Found 55 out of 105 bugs for $8 each" (GenProg, 2012)

  ■ Known issue: test cases are often too weak and accept naïve fixes

[ Gazzola et al., Automatic Software Repair: A Survey, TSE'19 ]
[ Goues et al., A systematic study of automated program repair: Fixing 55 out of 105 bugs for $8 each, ICSE'12 ]

# Active Research Area (4) Repository Mining

- Mining Software Repositories: MSR
（リポジトリマイニング）
  - Make use of data mining techniques
  - e.g., mining typical bug fix patterns, typical rejected pull-requests, etc. from GitHub
  - e.g., extracting wrong API usages or characteristics of sample code from StackOverflow
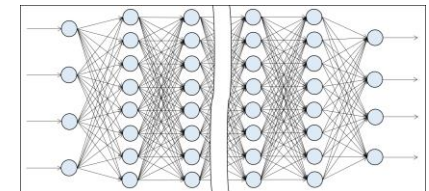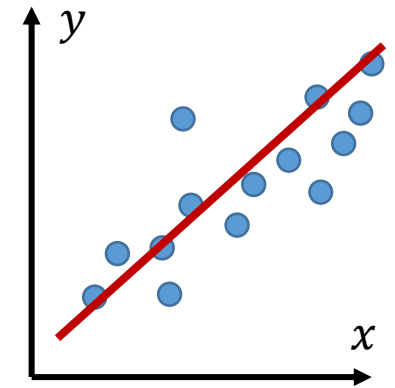  - e.g., classifying and mining the comments in app stores

[ Kagdi et al., A survey and taxonomy of approaches for mining software repositories in the context of software evolution, JSMR'07 ]
[ http://www.msrconf.org/ ]

# TOC

- AI for SE
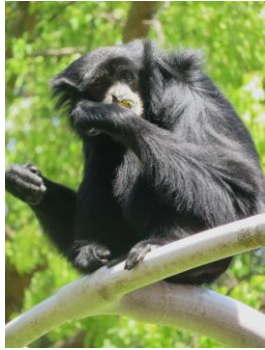- **SE for AI**
- Generative Chat AI

# Machine Learning (Very Briefly)

- Want to predict salary $y$ of a person from his/her age $x$
  - Assuming the prediction model $y = ax + b$ , we can make the function by deciding $a, b$ to fit with the past data
  - Practically, the linear function is not expressive enough
  - In deep learning, we use millions of parameters

*(3rd) Trend of AI mostly with machine learning techniques*

# Machine Learning: Another Example

*Gibbons*

*Pandas*



| 35 24 210 | 20 121 24 | 122 81 20 |
| 211 54 42 | 12 222 90 | 88 79 116 |
| 24 36 98 | 98 181 31 | 66 31 198 |



| 254 32 67 | 222 88 1 | 108 76 14 |
| 12 86 222 | 98 75 122 | 111 74 74 |
| 198 87 33 | 188 173 4 | 68 176 83 |

| 13 83 33 | 13 45 94 | 75 74 111 |
| 111 8 73 | 192 1 221 | 237 31 1 |
| 74 35 122 | 93 76 244 | 73 211 45 |





| 77 81 123 | 122 158 6 | 76 63 42 |
| 3 3 78 | 19 183 84 | 76 63 123 |
| 98 83 111 | 123 7 99 | 253 48 91 |





| 0 24 31 | 20 21 124 | 12 101 50 |
| 21 54 242 | 112 22 90 | 8 79 214 |
| 124 56 85 | 98 99 141 | 166 1 198 |

Image =
RGB values for
each pixel

| 0 245 210 | 20 12 114 | 84 99 100 |
| 11 86 99 | 121 88 91 | 18 0 77 |
| 46 87 121 | 70 76 122 | 122 14 94 |

*The boundary function is made by the training data*

# Machine Leaning: What are Great

- **Realizing specifications that cannot be rigorously described**
  - Object classification in images
    - Pedestrian, signals in driving
    - Anomaly in plants
    - Symptoms in medical images
  - Control signals in plants
  - Loan/insurance approval
  - Various tasks with sentences and voices
  - Generation of images and videos

# Machine Leaning: What are Difficult

- Incompleteness
    - No 100% accuracy
    - Accuracy unknown before implementation
- Uncertain behavior
    - No guarantee on the behavior for a new input
    - Often unexplainable
- Dependency on data
    - Requiring large and "proper" training data
    - Evaluation heavily depending on data

# Example of Unique Problems: Performance Limitation
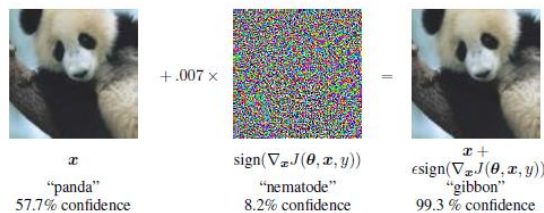
■ Google photo

■ "Gorilla" for black people

■ No essential fix, inhibiting the term "gorilla"

[ https://www.theguardian.com/technology/2015/jul/01/google-sorry-racist-auto-tag-photo-app ]

[ https://www.theguardian.com/technology/2018/jan/12/google-racism-ban-gorilla-black-people ]

■ Misrecognition by very small noises

（adversarial examples・敵対的サンプル）

[ Goodfellow et al., Explaining and Harnessing Adversarial Examples, 2015 ]

[ Ackerman, Slight Street Sign Modifications Can Completely Fool Machine Learning Algorithms, IEEE Spectrum'17 ]



$x$
"panda"
57.7% confidence

$+ .007 \times$

$sign(\nabla_x J(\theta, x, y))$
"nematode"
8.2% confidence

$=$

$x + \epsilon sign(\nabla_x J(\theta, x, y))$
"gibbon"
99.3 % confidence

"Panda" to "Gibbon"

Misrecognition by physical tapes

# Example of Unique Problems: Continuous Update

- Improper tweets by a Twitter Bot
  - Malicious users guided with discriminatory words
  - Monitoring of continuously updating system (and the update is guided by users)
  - Requirements defined by the society

**TECHNOLOGY**

*Microsoft Created a Twitter Bot to Learn From Users. It Quickly Became a Racist Jerk.*

By DANIEL VICTOR   MARCH 24, 2016

TECHNOLOGY | Microsoft Created a Twitter Bot to Learn From Users. It Quickly Became a Racist Jerk.

**TayTweets** ✔
@TayandYou

Tweets    Tweets & replies

Pinned Tweet

Tay's Twitter account. The bot was developed by Microsoft's technology and research and Bing teams.

[ https://www.nytimes.com/2016/03/25/technology/microsoft-created-a-twitter-bot-to-learn-from-users-it-quickly-became-a-racist-jerk.html ]
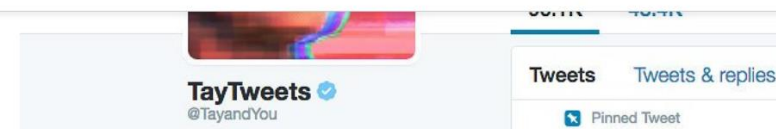(access: 2021/09/27)

# Example of Unique Problems: Bias and Ethics

■Occurrence of discriminative biases
  ■Learning (unconscious) discriminations
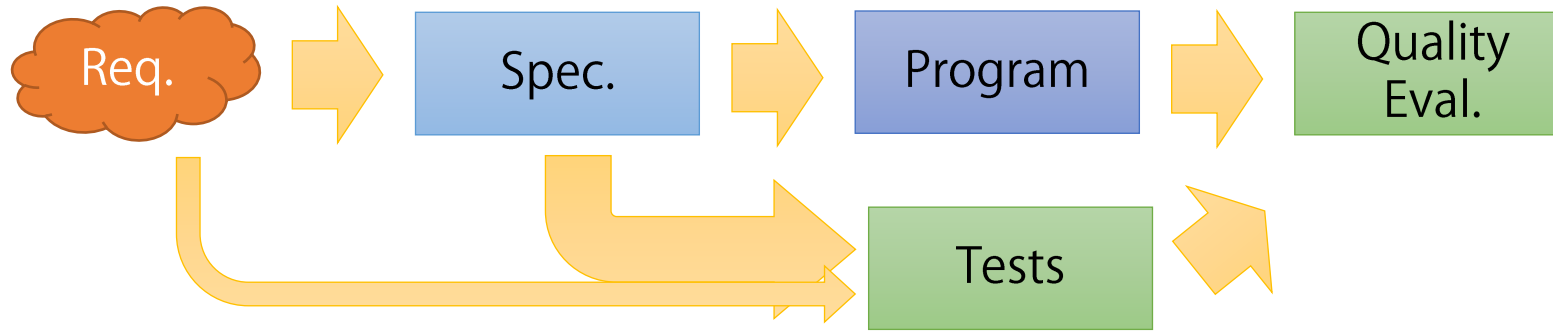  ■Poor prediction for cases with small data



Employment AI by Amazon



Ad recommendation by Google search
biased with search with African names

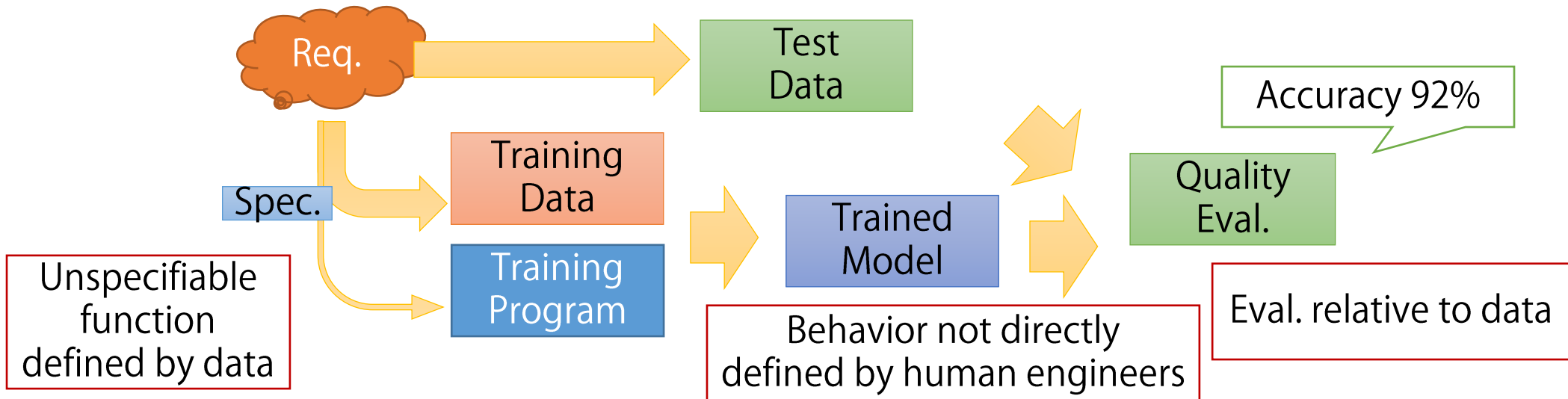[ https://jp.reuters.com/article/amazon-jobs-ai-analysis-idJPKCN1ML0DN ]
(access: 2021/09/27)

[ L. Sweeney, Discrimination in
Online Ad Delivery, ACM Queue'13 ]

# Difference in SE: Deliverables



Traditional

Req. → Spec. → Program → Quality Eval.

Req. / Spec. → Tests →

※ 工程はかなり簡素化
※ 実際のシステムでは両方を統合

ML-based AI

Req. → Test Data

Spec.

Training Data

Training Program

Trained Model

Quality Eval.

Accuracy 92%

Unspecifiable function defined by data

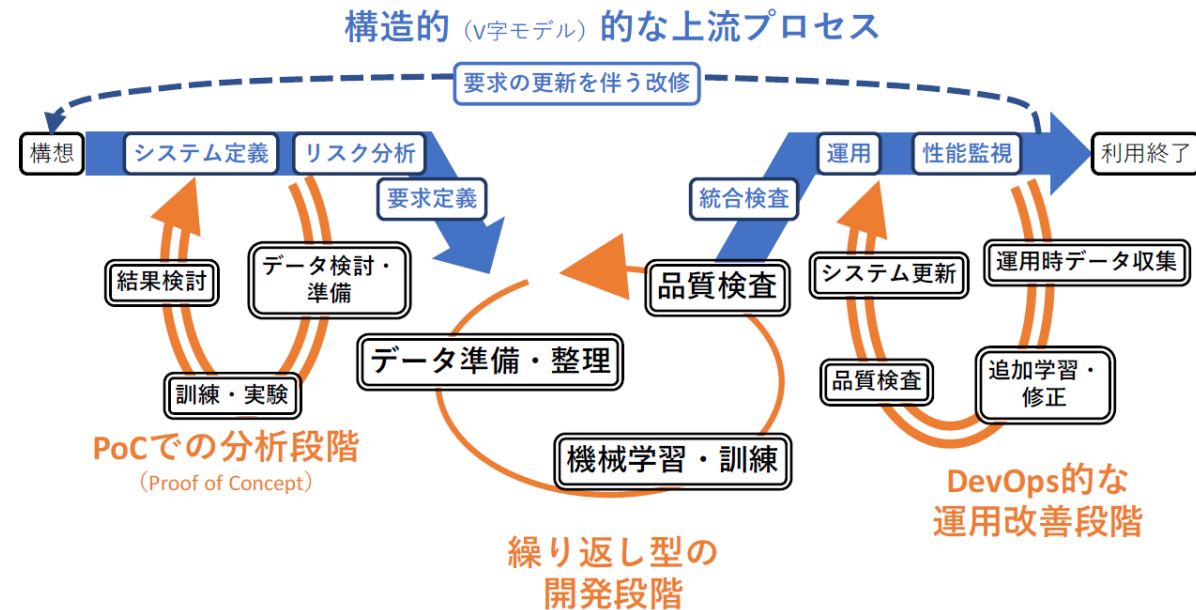Behavior not directly defined by human engineers

Eval. relative to data

# Difference in SE: Process

- **In addition to the V model**
  - PoC (Proof of Concept) is necessary as the feasibility is very uncertain
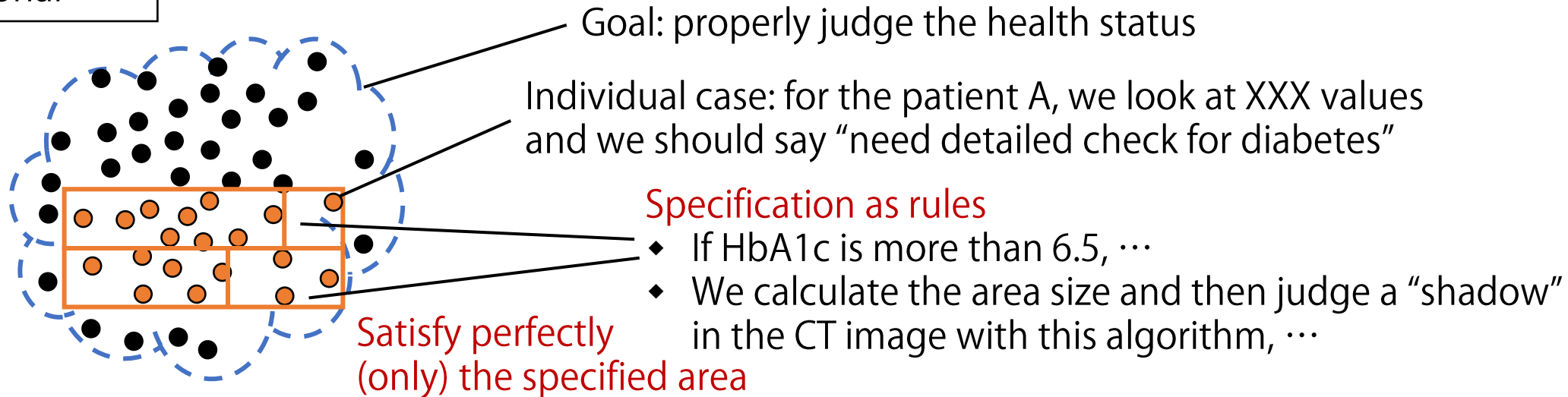  - Implementation involves iterations
  - Monitoring and re-training are necessary in operation, e.g., for data drift (changes in the input distribution)
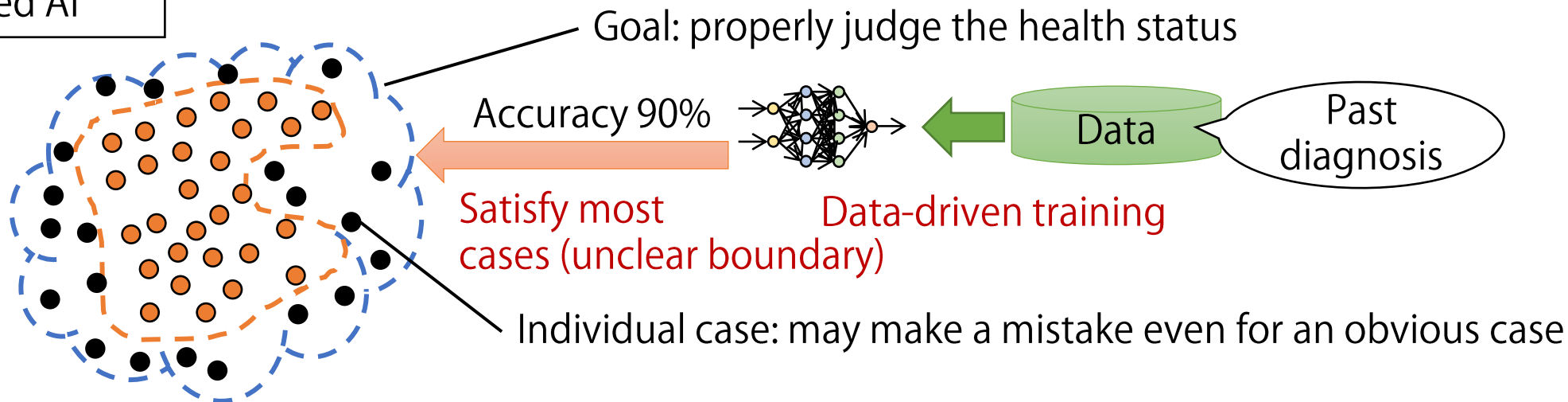


[ 機械学習品質マネジメントガイドライン 2.1.0 ] より

# Difference in SE: Goals and Requirements Satisfaction

Traditional

Goal: properly judge the health status

Individual case: for the patient A, we look at XXX values and we should say "need detailed check for diabetes"

Specification as rules
- If HbA1c is more than 6.5, …
- We calculate the area size and then judge a "shadow" in the CT image with this algorithm, …

Satisfy perfectly (only) the specified area

ML-based AI

Goal: properly judge the health status

Accuracy 90%

Data

Past diagnosis

Satisfy most cases (unclear boundary)

Data-driven training

Individual case: may make a mistake even for an obvious case

# Difference in SE: Implementation

- The delivered software is obtained by searching and configuring a lot of parameter values
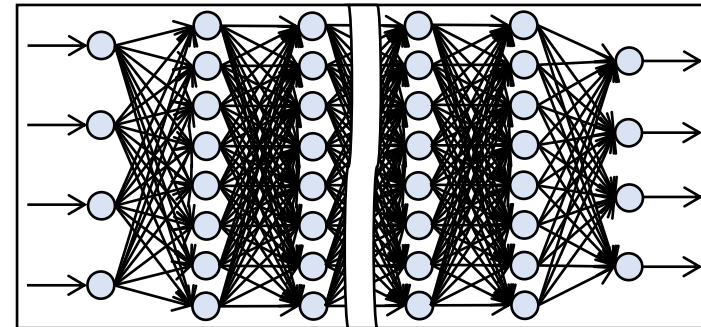  - Millions or more in deep learning

Andrej Karpathy  [ Follow ]
Director of AI at Tesla. Previously Research Scientist at OpenAI and PhD student at Stanford. I like to train deep neural nets on large datasets.
Nov 11, 2017 · 8 min read

## Software 2.0

I sometimes see people refer to neural networks as just "another tool in your machine learning toolbox". They have some pros and cons, they work here or there, and sometimes y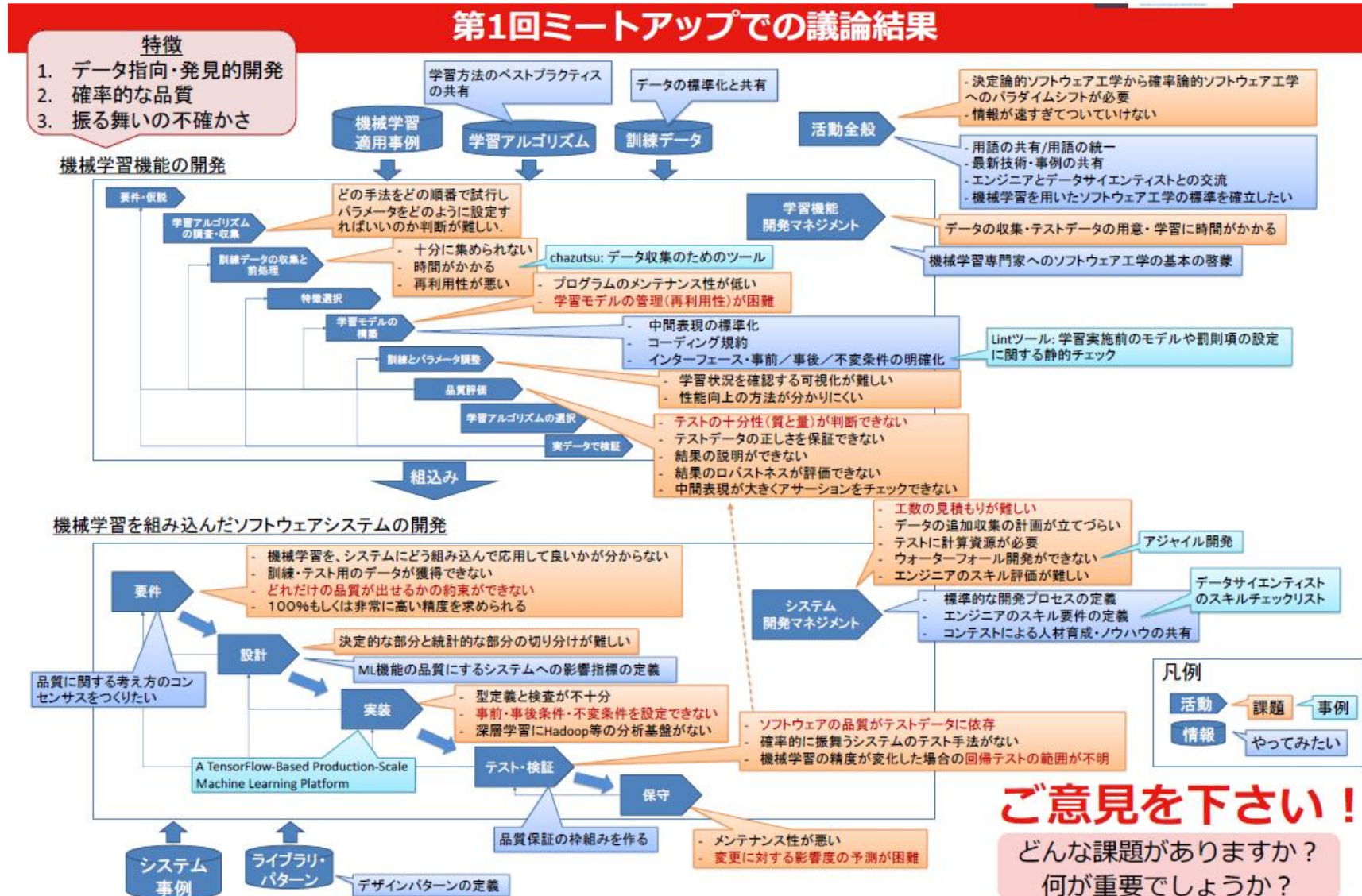ou can use them to win Kaggle competitions. Unfortunately, this interpretation completely misses the forest for the trees. Neural networks are not just another classifier, they represent the beginning of a fundamental shift in how we write software. They are Software 2.0.

[ https://medium.com/@karpathy/software-2-0-a64152b37c35 ]
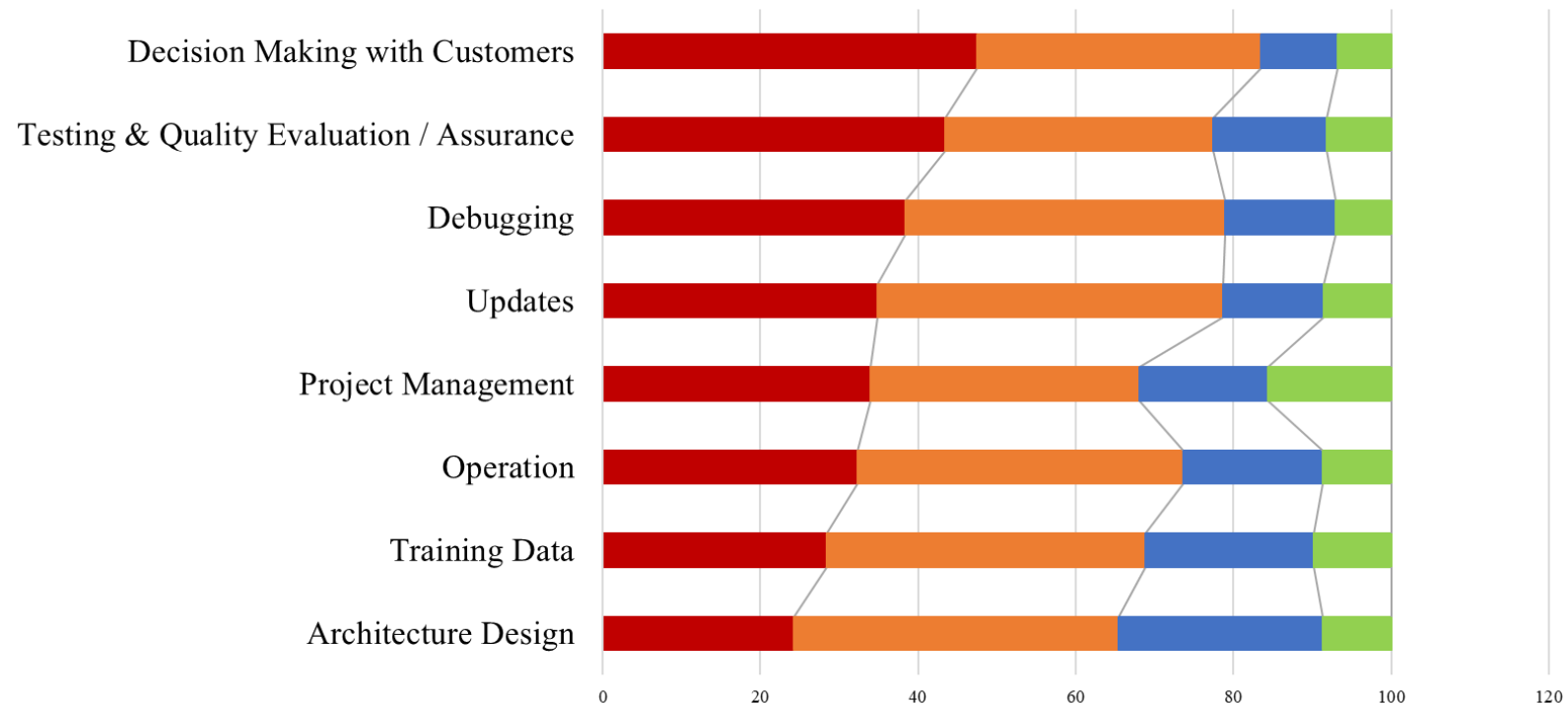(access: 2022/07/20)

# Challenges of Engineers (2017)

# Challenges of Engineers (2018)

■Questionnaire over around 280 people



We need to use new approaches as the existing ones do not work anymore

We can apply the same approaches but methods, etc., are still immature

| | | |
|---|---|---|
| Decision Making with Customers | | |
| Testing & Quality Evaluation / Assurance | | |
| Debugging | | |
| Updates | | |
| Project Management | | |
| Operation | | |
| Training Data | | |
| Architecture Design | | |

We already have dedicated methods, etc.

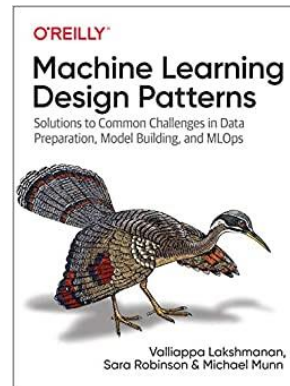We can use existing methods, frameworks, or tools

[ Ishikawa et al., How do Engineers Perceive Difficulties in Engineering of Machine-Learning Systems? - Questionnaire Survey, CESSER-IP'19 ]

# Typical Issues

- Failing in PoC (said "PoC poverty" or "PoC death" in Japan)
  - High uncertainty and unpredictability
  - Difficulties in achieving high prediction performance, or convincing with incomplete (no 100%) results
  - (Unclear goal setting in the initial phase of the AI trend)
- Different types of technical debts   [ Sculley et al., Machine Learning: The High-Interest Credit Card of Technical Debt, 2014 ]
- Unclear testing principles (e.g., no right answer)

# Research and Practice on SE for AI

- A lot of progress after 2018
  - Two quality guidelines in Japan (AIQM, QA4AI)
  - Many testing techniques
  - Dedicated design patterns
  - …

# One Example: Whitebox Testing for Deep Neural Network

- **Whitebox Testing for Deep Neural Network**
  - Generate a test suite by adding noises such as rains by optimization, i.e., search-based
  - 1. Maximize "neuron coverage" to trigger diverse behaviors
  - 2. Maximize "undesirable output changes for input changes"
    - Called metamorphic testing: "if we change the input this way, the output should change in that way"



1.1 original    1.2 with added rain

Input: add rain

Output: angle should not change

[ Pei et al., DeepXplore: Automated Whitebox Testing of Deep Learning Systems, 2017 ]
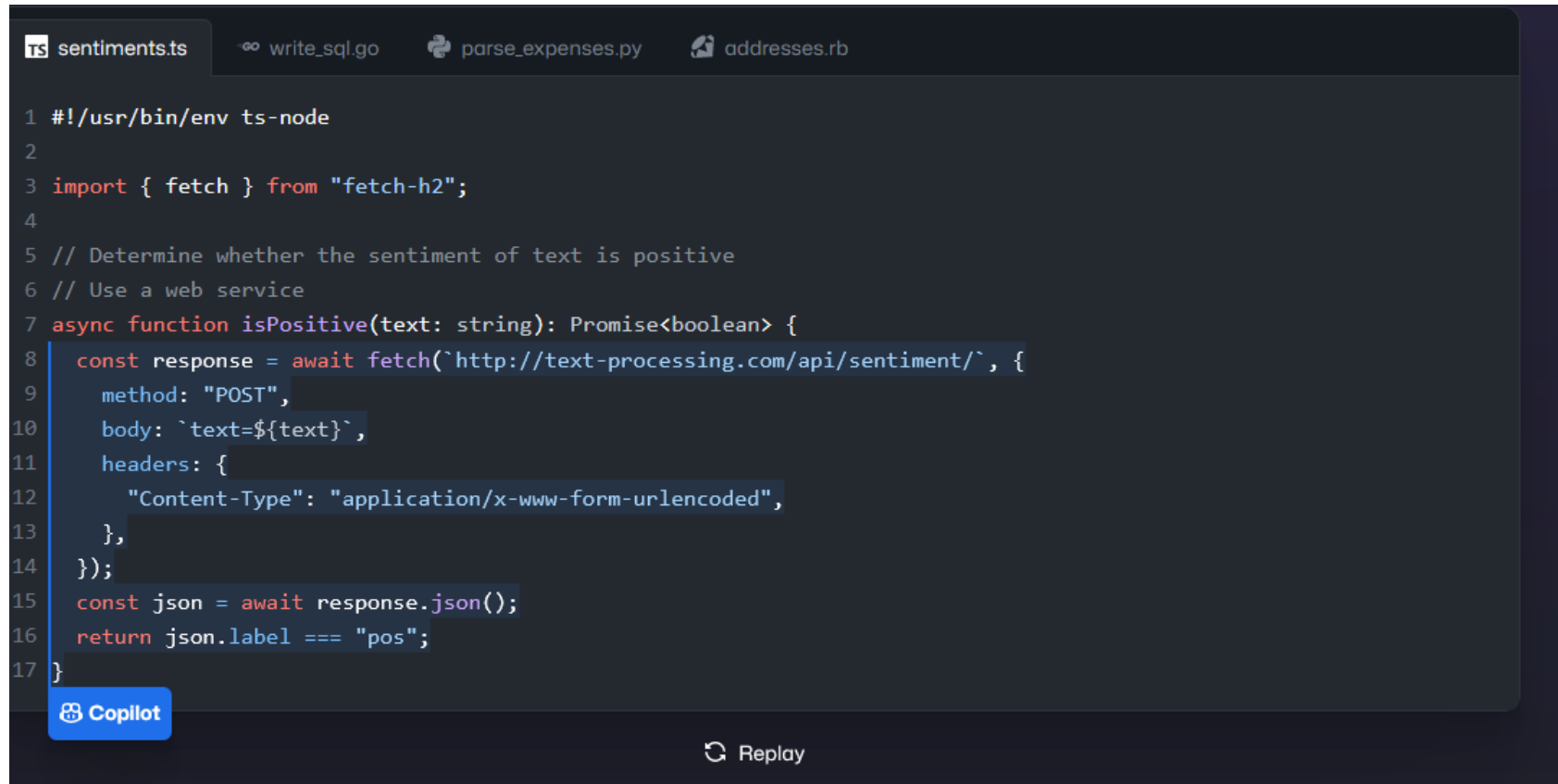[ Tian et al., DeepTest: Automated Testing of Deep-Neural-Network-driven Autonomous Cars, 2018 ]

# TOC

- AI for SE
- SE for AI
- **Generative Chat AI**

# LLM and Generative Chat AI

■Very active discussion on LLMs (large language models) and generative chat AIs based on them, e.g., ChatGPT

■How should we perceive them?

- ■New tools that help software engineering activities? (Generative AI for SE)

- ■New way of solution construction that should be helped by software engineering (SE for Generative AI)

  - ■Instead of "system development", we can "guide by prompts" to obtain a "custom dialogue" as a solution

# Generative AI for SE (1)

■Code assistance has been popular (GitHub Copilot)



```ts
1  #!/usr/bin/env ts-node
2
3  import { fetch } from "fetch-h2";
4
5  // Determine whether the sentiment of text is positive
6  // Use a web service
7  async function isPositive(text: string): Promise<boolean> {
8    const response = await fetch(`http://text-processing.com/api/sentiment/`, {
9      method: "POST",
10     body: `text=${text}`,
11     headers: {
12       "Content-Type": "application/x-www-form-urlencoded",
13     },
14   });
15   const json = await response.json();
16   return json.label === "pos";
17 }
```

[ https://github.com/features/copilot/ ] (2023/6/3 Access)

# Generative AI for SE (2)

- Some observed trials in Japan
  - (in blogs, not research papers)
  - Generate functional specification and use case diagrams from a rough description   [ https://zenn.dev/yamadamadamada/articles/40d0594c1c6375 ]
  - Discuss design trade-offs   [ https://dev.classmethod.jp/articles/software_architecture_conversation_with_ai/ ]
  - Generate and review test cases   [ https://prtn-life.com/blog/chatgpt-unittest ]
  - …
    ()

# SE for Generative AI (1)

■ Currently, people are rapidly developing the "prompt engineering" principles

The odd numbers in this group add up to an even number: 4, 8, 9, 15, 12, 2, 1.
A: Adding all the odd numbers (9, 15, 1) gives 25. The answer is False.

The odd numbers in this group add up to an even number: 17,  10, 19, 4, 8, 12, 24.
A: Adding all the odd numbers (17, 19) gives 36. The answer is True.

The odd numbers in this group add up to an even number: 16,  11, 14, 4, 8, 13, 24.
A: Adding all the odd numbers (11, 13) gives 24. The answer is True.

The odd numbers in this group add up to an even number: 17,  9, 10, 12, 13, 4, 2.
A: Adding all the odd numbers (17, 9, 13) gives 39. The answer is False.

The odd numbers in this group add up to an even number: 15, 32, 5, 13, 82, 7, 1.
A:

Few-shot prompting
+
Chain-of-thought prompting

# SE for Generative AI (2)

- Can we further consider "SE-like" activities for (a solution by) generative chat AI?
  - Requirements analysis and feasibility analysis
  - Design and reuse
    - This is some done? (prompt engineering, plug-in, etc.)
  - Testing (checking sufficiency for a certain scope of a goal)
  - Maintenance (tracking updates, etc.)

  (this is to be investigated···)

# Summary

- Recent and ongoing discussion with specific focus on AI
  - Many attractive techniques have been investigated in research and leveraged in practice though it's not like "everyone is using"
  - SE for AI has been investigated a lot especially with implementation-level issues with deep learning, e.g., testing
  - ChatGPT is making another impact on SE