# Beyond Code Generation
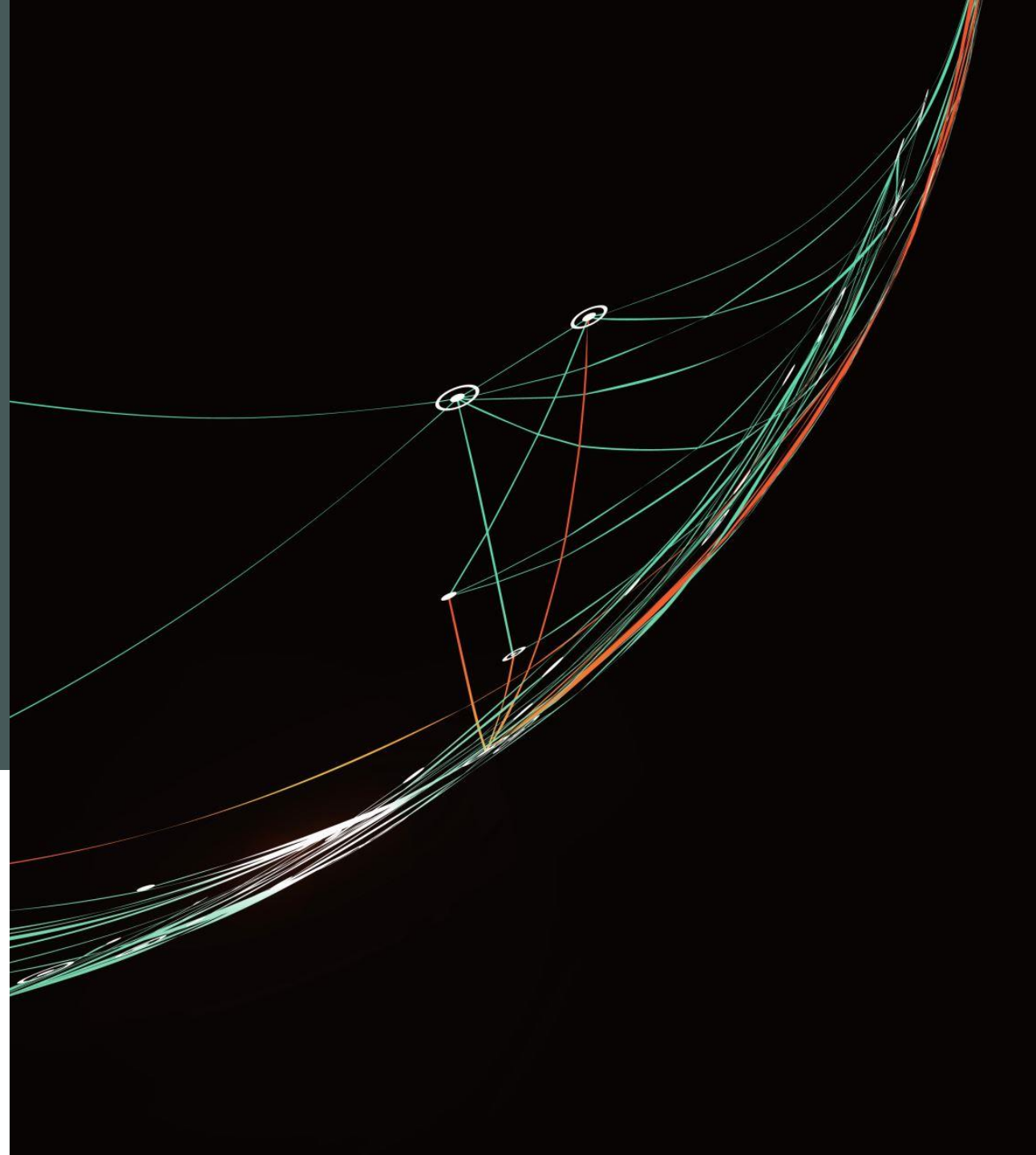
# Towards Next-Generation AI for SE
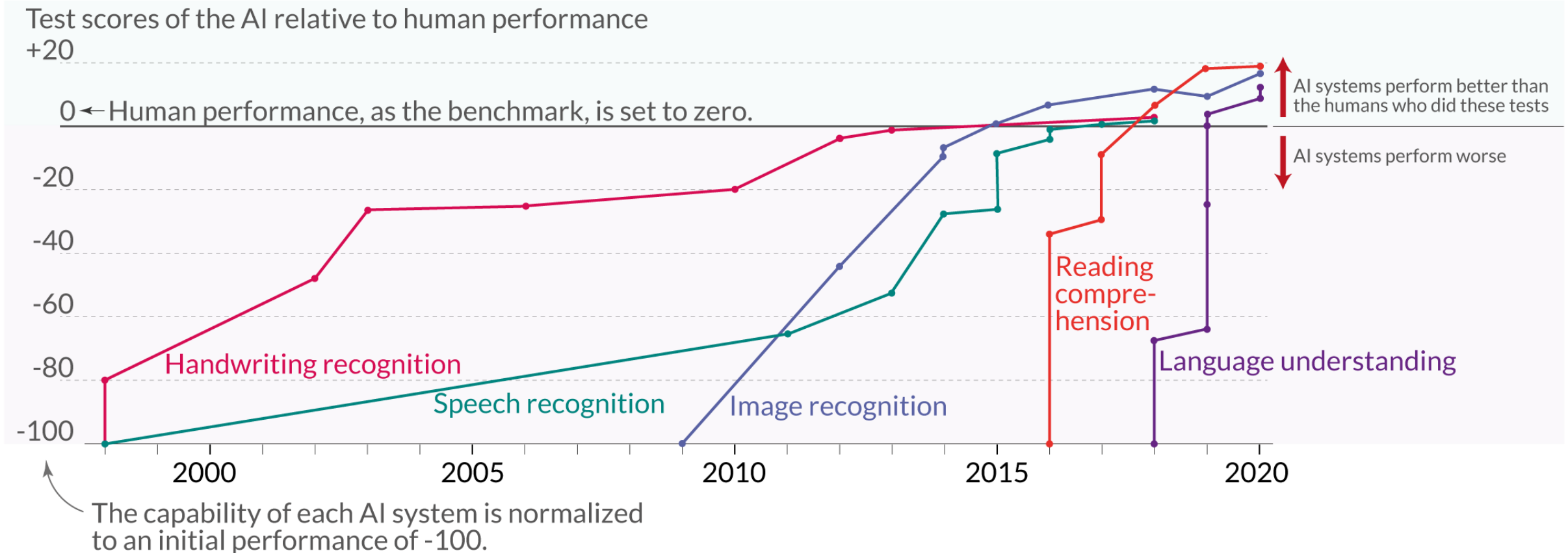
Vincent J. Hellendoorn

December 3rd, 2023
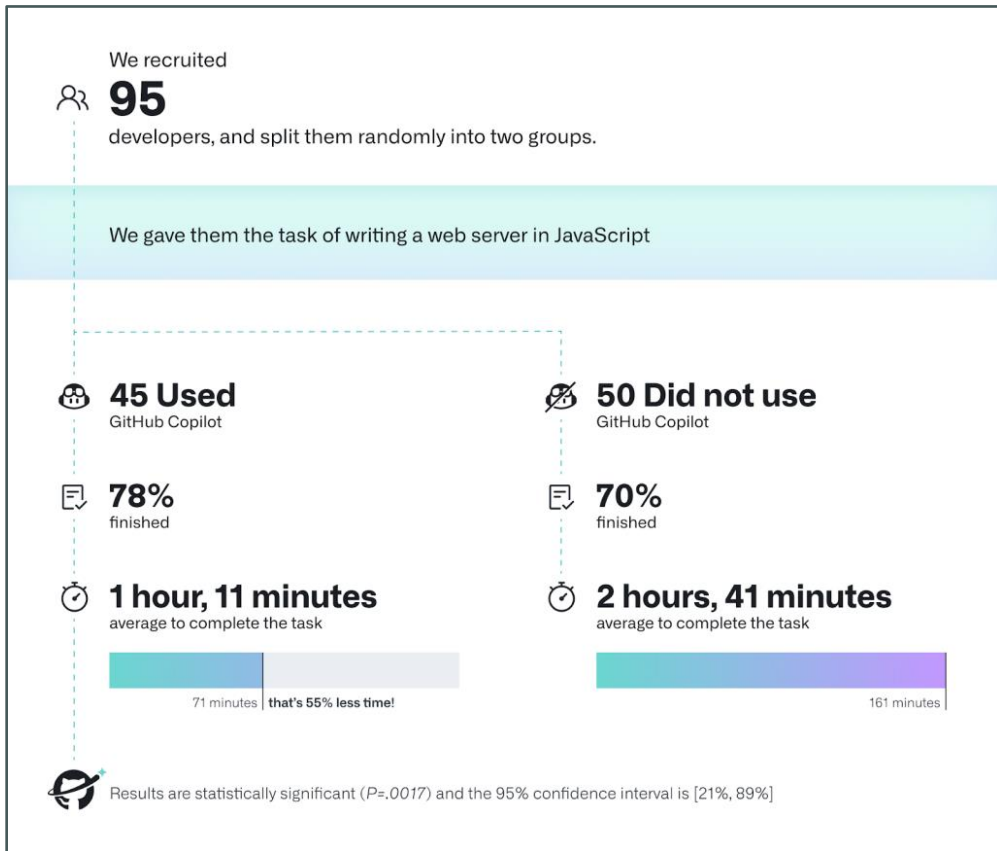
MAPS @ FSE

# Where We Were



Test scores of the AI relative to human performance

AI systems perform better than the humans who did these tests

AI systems perform worse

0 ← Human performance, as the benchmark, is set to zero.

Handwriting recognition

Speech recognition

Image recognition

Reading compre-hension

Language understanding

The capability of each AI system is normalized to an initial performance of -100.

# Where We Were



We recruited

## 95
developers, and split them randomly into two groups.

We gave them the task of writing a web server in JavaScript

**45 Used**
GitHub Copilot

**50 Did not use**
GitHub Copilot

**78%**
finished

**70%**
finished

**1 hour, 11 minutes**
average to complete the task

**2 hours, 41 minutes**
average to complete the task

71 minutes | that's 55% less time!

161 minutes

Results are statistically significant ($P=.0017$) and the 95% confidence interval is [21%, 89%]



```go
1  package main
2
3  type CategorySummary struct {
4      Title        string
5      Tasks        int
6      AvgValue     float64
7  }
8
9  func createTables(db *sql.DB) {
10     db.Exec("CREATE TABLE tasks (id INTEGER PRIMARY KEY, title TEXT, value INTEGER, category TEXT
11 }
12
13 func createCategorySummaries(db *sql.D
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
```

# Where We Aren't

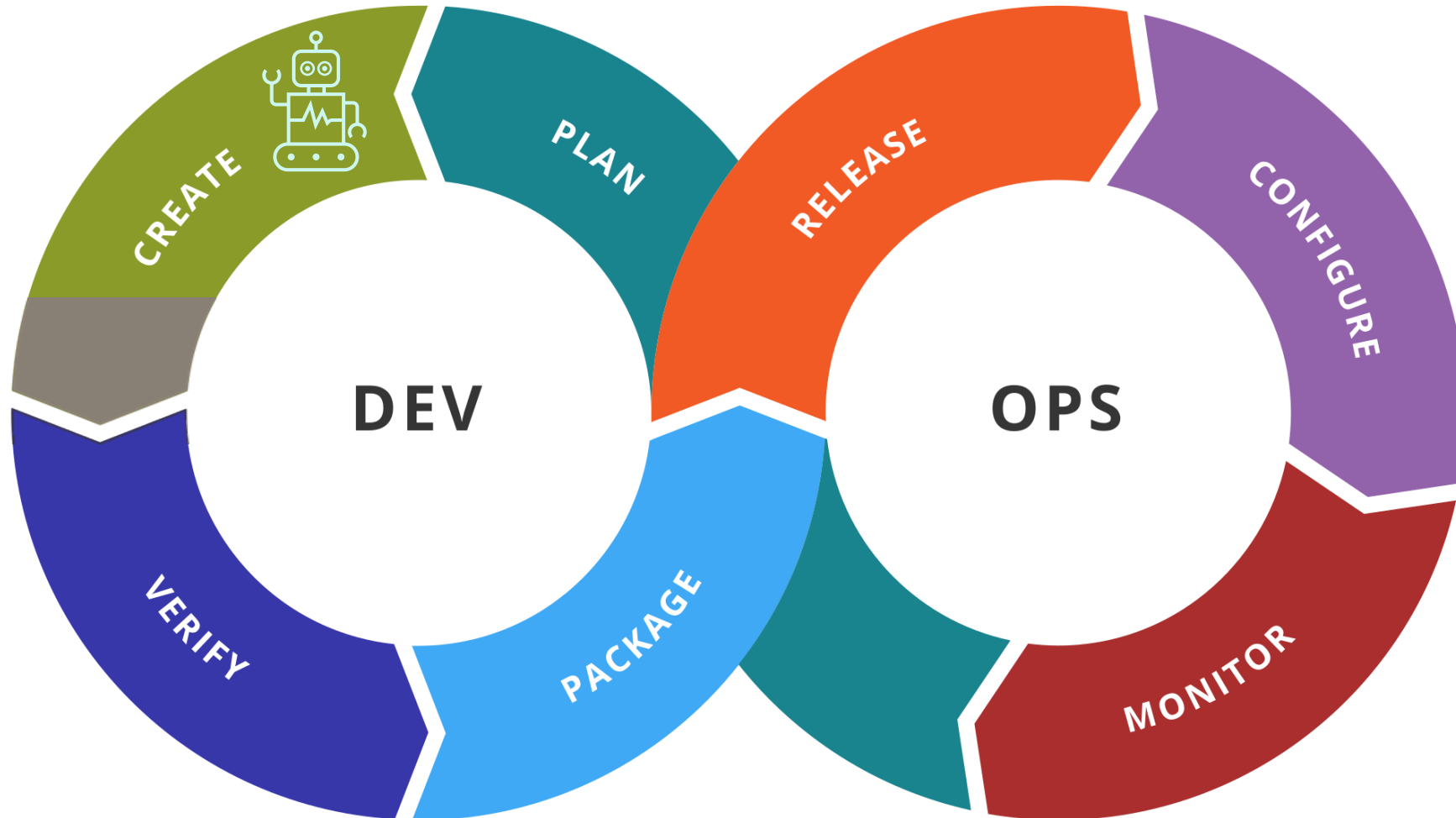How fairly rated is AI?

Fairly rated — 23.4%

Underrated — 25.1%

Overrated — 51.6%

# Today:
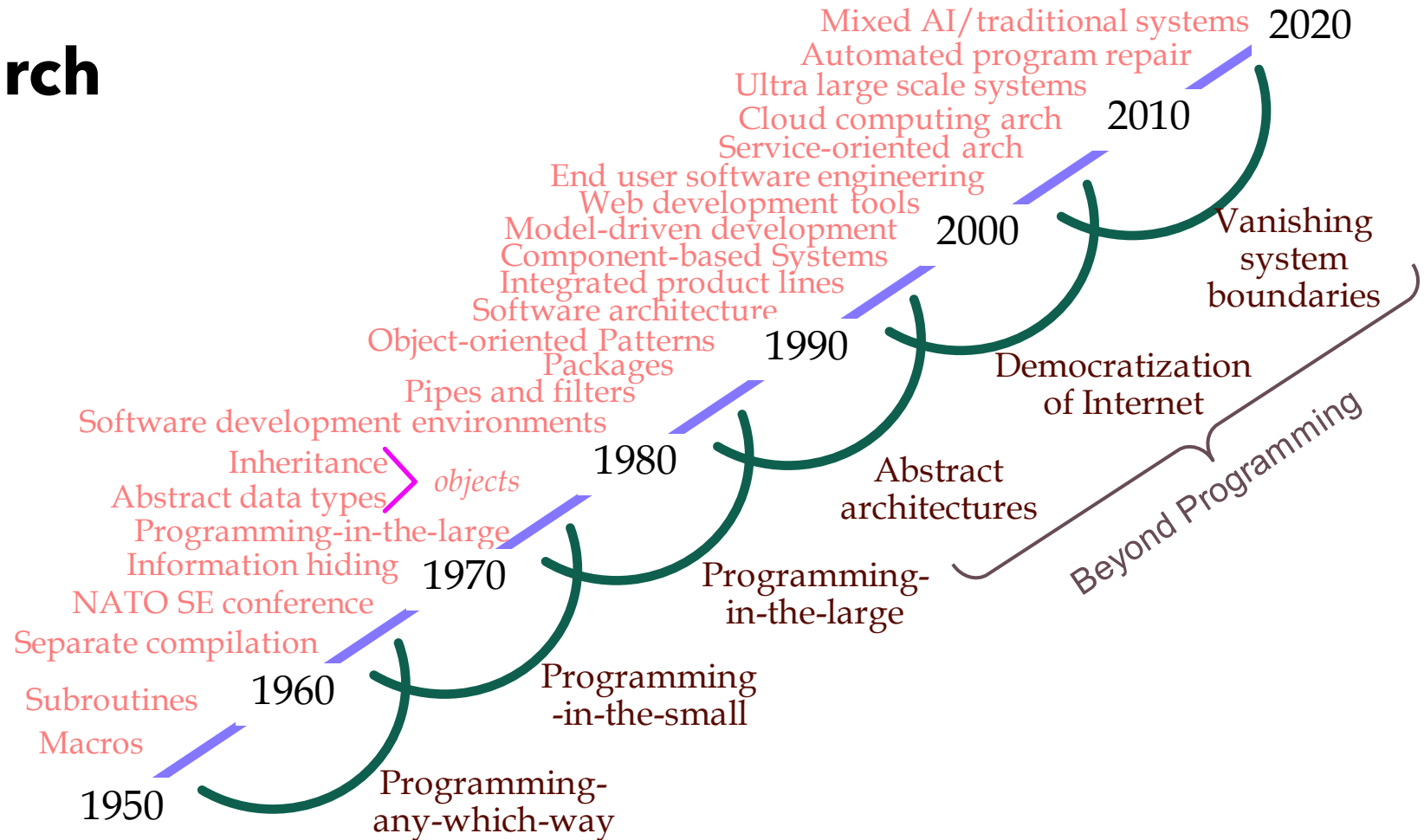# Where We Are, and Are Going

# It's time to move beyond writing assistants

1. Start by reflecting on the olden days (<2020)

2. Then, discuss how LLMs changed the picture

3. Next, highlight recent progress & trends

4. Finally, promises, challenges, needs & tips

# Software Development Always Changes

And so does **SE research**

Mixed AI/traditional systems — 2020
Automated program repair
Ultra large scale systems
Cloud computing arch — 2010
Service-oriented arch
End user software engineering
Web development tools
Model-driven development — 2000
Component-based Systems
Integrated product lines
Software architecture
Object-oriented Patterns — 1990
Packages
Pipes and filters
Software development environments
Inheritance — 1980
Abstract data types ⟩ *objects*
Programming-in-the-large
Information hiding — 1970
NATO SE conference
Separate compilation
Subroutines — 1960
Macros
1950

Vanishing system boundaries

Democratization of Internet

Abstract architectures

Programming-in-the-large

Programming-in-the-small

Programming-any-which-way

Beyond Programming

# "First Wave" of ML for SE

**Hand-extracted features** fed to **off-the-shelf learners**

Multiple, generic models (e.g., decision trees)

Practitioner focuses on **features & statistics**

# "First Wave" of ML for SE

What that looked like:

| Name | Description |
|------|-------------|
| NR | Number of revisions |
| NREF | Number of times a file has been refactored |
| NFIX | Number of times a file was involved in bug-fixing |
| NAUTH | Number of authors who committed the file |
| LINES | Lines added and removed (sum, max, average) |
| CHURN | Codechurn (sum, maximum and average) Codechurn is computed as $\sum_R(addedLOC - deletedLOC)$, where $R$ is the set of all revisions |
| CHGSET | Change set size, *i.e.,* number of files committed together to the repository (maximum and average) |
| AGE | Age (in number of weeks) and weighted age computed as $\frac{\sum_{i=1}^N Age(i) \times addedLOC(i)}{\sum_{i=1}^N addedLOC(i)}$, where $Age(i)$ is the number of weeks starting from the release date for revision $i$, and $addedLOC(i)$ is the number of lines of code added at revision $i$ |

| Category of approach | GLM | | DT | | NB | |
|----------------------|-----|-----|-----|-----|-----|-----|
| | Mean | Var | Mean | Var | Mean | Var |
| Process metrics (MOSER) | 6.4 | 0.64 | 6.2 | 1.36 | 7.2 | 3.44 |
| Previous defects (BUG-CAT) | 6.6 | 5.84 | 4.2 | 3.76 | 5.2 | 10.16 |
| Entropy of changes (HCM, WHCM, EDHCM, LDHCM, LGDHCM) | 5.8 | 12.16 | 4.6 | 6.64 | 6.8 | 10.56 |
| Code metrics (CK+OO) | 9.4 | 0.64 | 7.4 | 5.84 | 9.2 | 1.36 |
| Churn of code metrics (LGDCHU) | 8.8 | 0.96 | 5.6 | 5.84 | 6.4 | 5.84 |
| Entropy of code metrics (LDHH) | 9.0 | 0.81 | 6.6 | 1.84 | 7.2 | 8.96 |

D'ambros, ESE, Evaluating Defect Prediction Approaches: A Benchmark and an Extensive Comparison

# "First Wave" of ML for SE

**Hand-extracted features** fed to **off-the-shelf learners**

**Pros:**

- Useful for **almost any task** where decisions are made

**Cons:**

- Feature selection **limits performance**, requires manual effort

- Largely **inapplicable to code**

# "Second Wave" of ML for SE
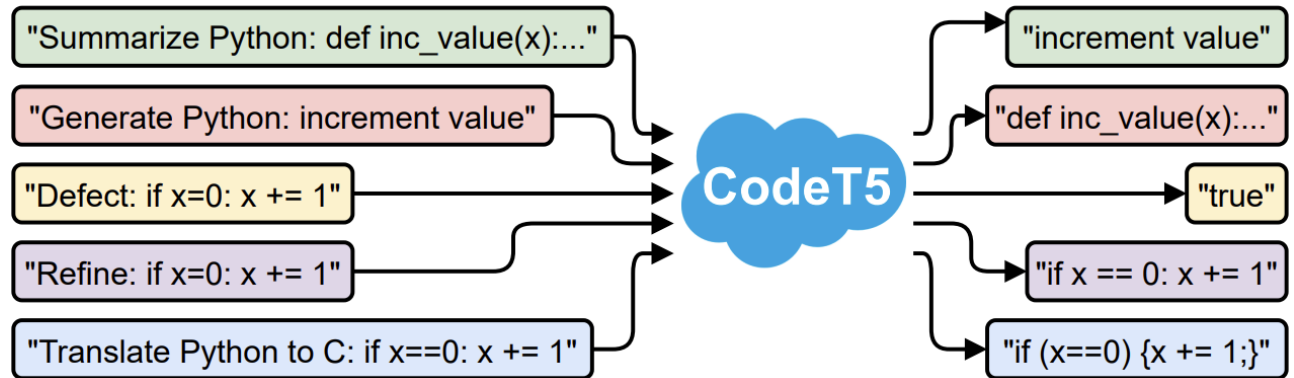
## Learning **from** and for **Source Code**

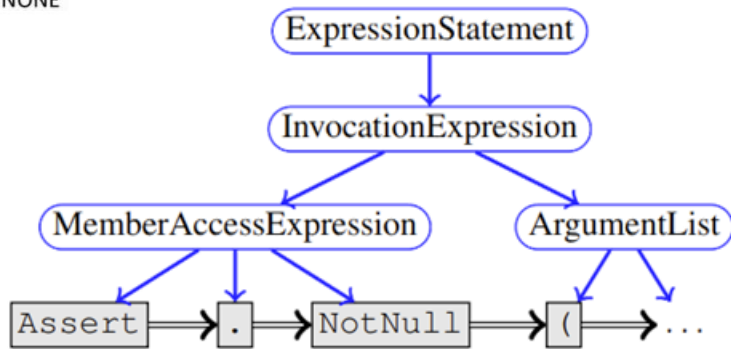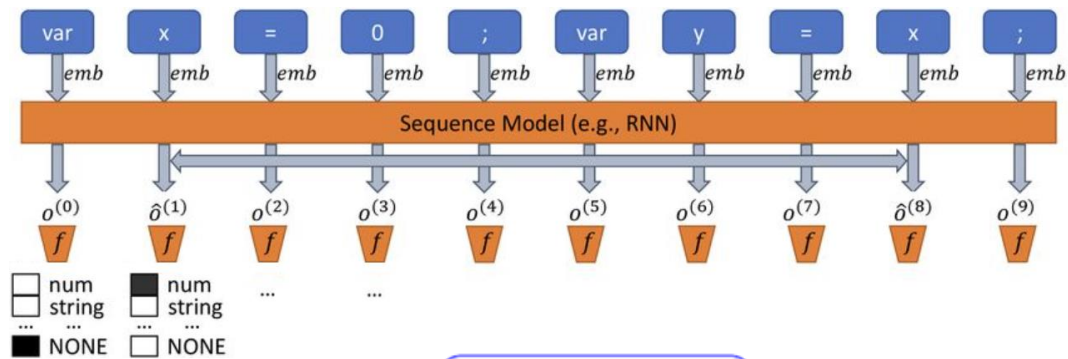Feature learning is left to the models

Model design is often **inspired by tasks**

Image by Stable Diffusion XL

# "Second Wave" of ML for SE

## What that looked like:

Marc Brockschmidth, MSR. (Deep) Learning from Programs. Slides from: https://slideplayer.com/slide/16532816/
Allamanis et al., ICLR'18. Learning to Represent Programs with Graphs. https://arxiv.org/pdf/1711.00740.pdf
Wang et al., EMNLP'21. CodeT5: Identifier-aware Unified Pre-trained Encoder-Decoder Models for Code Understanding and Generation

# "Second Wave" of ML for SE

Learning **from** and for **Source Code**

**Pros:**

• Increased **expressivity**, less manual effort

**Cons:**

• Learning limited by **dataset size**

• Limited **practical utility**

Image by Stable Diffusion XL

# Obviously, the "third wave" is learning program semantics, right?

Right?



Next Generation Models of Code

Code is static <u>and</u> dynamic
• Experts think of both
• Our models should too



## 6.1 The Third Wave of Machine Learning

The first wave of machine learning for source code applied off-the-shelf machine learning tools with hand-extracted features. The second wave, reviewed here, avoids manual feature extraction and uses the source code itself within machine learning heavily drawing inspiration from existing machine learning methods in NLP and elsewhere. The third wave promises new machine learning models informed by programming language semantics. What form will it take?

# The Bitter Lesson Strikes Again

GPT-3 (Brown et al., 2020), https://arxiv.org/pdf/2005.14165.pdf

# "Third Wave" of ML for SE

## Pretraining at **immense scale**

Not just on code

# "Third Wave" of ML for SE

## Pretraining at **immense scale**

Not just on code

In the world of **Large Language Models**, the goal is generation and the currency is compute



We trained this one early on

# "Third Wave" of ML for SE

Pretraining at **immense scale**

**Pros:**

• Can generate **large volumes of code** & text

• Extraordinary **representational power**

**Cons:**

• Inherently **generative**

• Extremely **data, compute hungry**

Image by Stable Diffusion XL

# Well, We Kept Feeding It



Artificial intelligence: Performance on knowledge tests vs. training computation

Performance on knowledge tests is measured with the MMLU benchmark[1]. Training computation is measured in total petaFLOP, which is $10^{15}$ floating-point operations[2].

# Where We Are: an Inflection Point



We used to model **any process**, but the models were bad

Now we can model **one process insanely well**

Let's spread the love

# Five Challenges for "AI for Code"

# Five Challenges for "AI for Code"

1. Expanding Context

# Five Challenges for "AI for Code"

1. Expanding Context

2. Software Maintenance

# Five Challenges for "AI for Code"

1. Expanding Context
2. Software Maintenance
3. Modeling Semantics

# Five Challenges for "AI for Code"

1. Expanding Context

2. Software Maintenance

3. Modeling Semantics

4. Interacting in Teams

# Five Challenges for "AI for Code"

1. **Expanding Context**

2. Software Maintenance

3. Modeling Semantics

4. Interacting in Teams

5. Navigating Process

# Context Scaling



LLM contexts have **grown rapidly**

Isn't attention cost **quadratic in input length**?

A: Who cares

- It's a **small fraction** of overall compute

- Often **worth the cost** during inference

https://cobusgreyling.medium.com/rag-llm-context-size-6728a2f44beb

# Context Scaling

## Context is what you **make of it**

- Pretraining "packs" multiples files together

- Shuffle GitHub? The model learns to ignore adjacent files


code_foo.py


code_bar.py


test_foo.py


test_bar.py

# Context Scaling

Context is what you **make of it**

- Pretraining "packs" multiples files together

- Shuffle GitHub? The model learns to ignore adjacent files

- Which is a waste

code_foo.py

code_bar.py

test_foo.py

test_bar.py

Rao et al., ASE'23. CAT-LM 🐱 Training Language Models on Aligned Code And Tests

# Using Long Contexts: Test Prediction

## Finding relevant data

Only 1M code/test pairs on GitHub?

No problem! Just **use everything**



Rao et al., ASE'23. CAT-LM 🐱 Training Language Models on Aligned Code And Tests

# Using Long Contexts: Test Prediction

## Finding relevant data

Some files are very long?

Throw **compute** at it!



Rao et al., ASE'23. CAT-LM 🐱 Training Language Models on Aligned Code And Tests

# Using Long Contexts:
# Test Prediction

## Of course this helps

We generate **more valid tests** on a fraction of the budget

- Coverage approaches, but **doesn't quite match**, human programmers

- Larger models could do better



Work led by Nikitha Rao (nikitharao@cmu.edu)
Cloud compute credits contributed by Google

Rao et al., ASE'23. CAT-LM 🐱 Training Language Models on Aligned Code And Tests

# Context Scaling – Lessons

Spend **compute** where it is due

Scaling **parameters** and **data** are slowing down. **Context** has a lot to offer.

Add **value through data**

Better a **million good tokens** than a trillion boring ones

Still, a **long way to go** for modeling real contexts

See SWEBench

# Five Challenges for "AI for Code"

1. Expanding Context
2. **Software Maintenance**
3. Modeling Semantics
4. Interacting in Teams
5. Navigating Process

# Software Maintenance

## Developers don't only Write Code

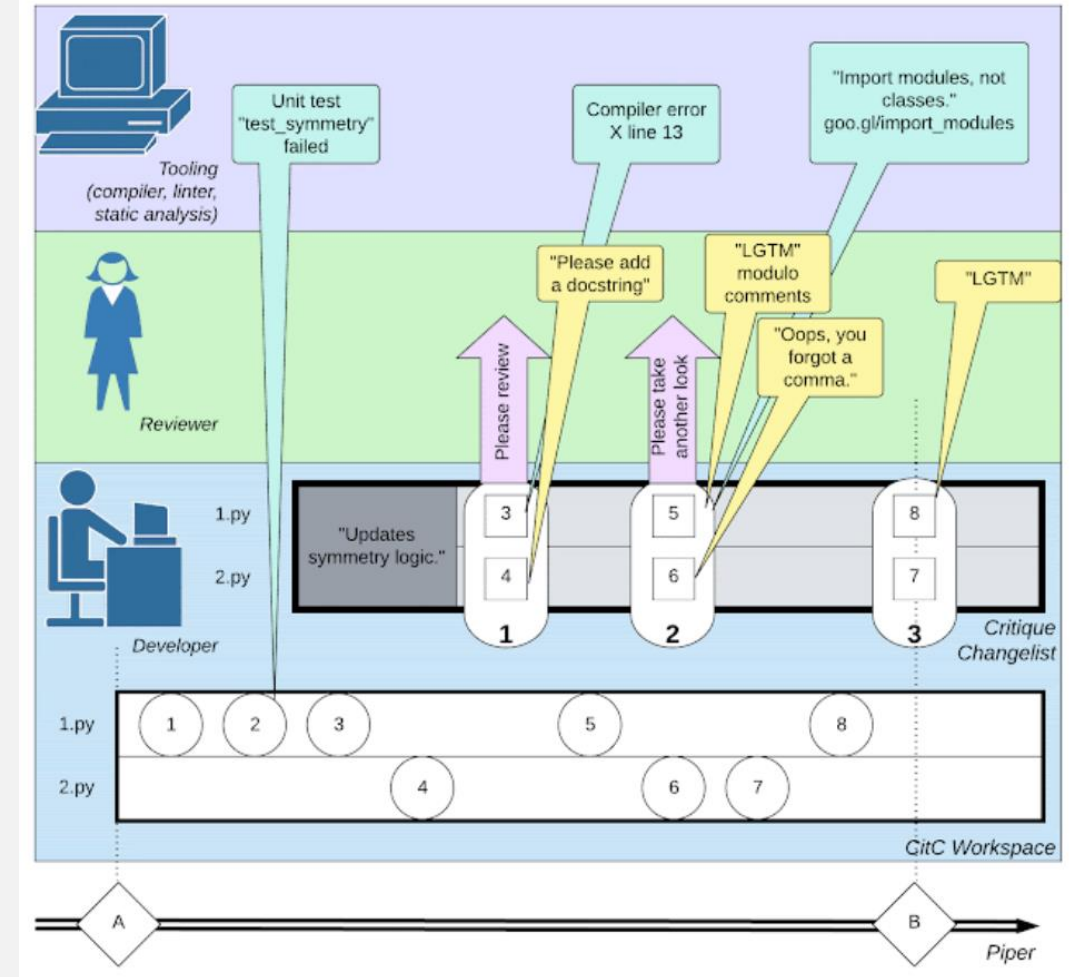Test prediction is still **generative**. What about **understanding, analyzing** code?

# Software Maintenance: Fault Localization

## A Case Study

All code has **bugs**, but most code isn't buggy.



Yang et al., ICSE'24. Large Language Models for Test-Free Fault Localization. https://arxiv.org/abs/2310.01726

# Software Maintenance: Fault Localization

## A Case Study

All code has **bugs**, but most code isn't buggy.

LLM have read all code. Can they tell?



Yang et al., ICSE'24. Large Language Models for Test-Free Fault Localization. https://arxiv.org/abs/2310.01726

# Software Maintenance: Fault Localization

## Generation ≠ Interpretation

But, the information **is already there**



**Causal Attention**

Left-to-right Large Language Model

Pretrained Model

LMFL

Yang et al., ICSE'24. Large Language Models for Test-Free Fault Localization. https://arxiv.org/abs/2310.01726

## Generation ≠ Interpretation

In this work, we train:

- Low-dimensional, bidirectional **adapter layers**

- A lightweight classifier



Yang et al., ICSE'24. Large Language Models for Test-Free Fault Localization. https://arxiv.org/abs/2310.01726

# Software Maintenance:
# Fault Localization

Improves over **purpose-built** methods



Work led by Aidan Yang (aidan@cmu.edu)

Yang et al., ICSE'24. Large Language Models for Test-Free Fault Localization. https://arxiv.org/abs/2310.01726

# Software Maintenance: Fault Localization

Improves over **purpose-built** methods

Drastically improves over the **initial LLM**



Work led by Aidan Yang (aidan@cmu.edu)

Yang et al., ICSE'24. Large Language Models for Test-Free Fault Localization. https://arxiv.org/abs/2310.01726

# A Theme Emerges

## Lean on the **pretrained model**

LLMs have seen **10,000x** what we can read in a lifetime

A ton of knowledge is **untapped** in those weights

To **tap in**, research needs to leverage on SE knowledge

# Five Challenges for "AI for Code"

1. Expanding Context
2. Software Maintenance
3. ~~Modeling Semantics~~
4. **Interacting in Teams**
5. Navigating Process

# Modeling the SE Process

Software development involves a lot **more than coding.**

Can we model the whole process?

# Modeling the SE Process

A more **holistic** view

Many tasks are **connected**

Developers appreciate help
anywhere, but **good UX is key**



**DIDACT**

Build Error Prediction

Build Repair

Comment Prediction

Tip Prediction

Comment Resolution

Edit Prediction

Variable Renaming

History-Augmented
Code Completion

# Modeling the SE Process

A few more examples of **promising trends**

# A Call to Action

The field needs **concerted efforts** to study AI in the wild

AI is inherently **about people** – observational studies are becoming key

LLMs have entered **supercomputing territory** – no lab can do this alone

Come together and build **high-quality benchmarks & tools**

# Five Challenges for "AI for Code"

1. Expanding Context
2. Software Maintenance
3. Modeling Semantics
4. Interacting in Teams
5. **Navigating Process**

# What's Next?

# What's Next? – Shifting Right



https://en.wikipedia.org/wiki/DevOps_toolchain#/media/File:Devops-toolchain.svg

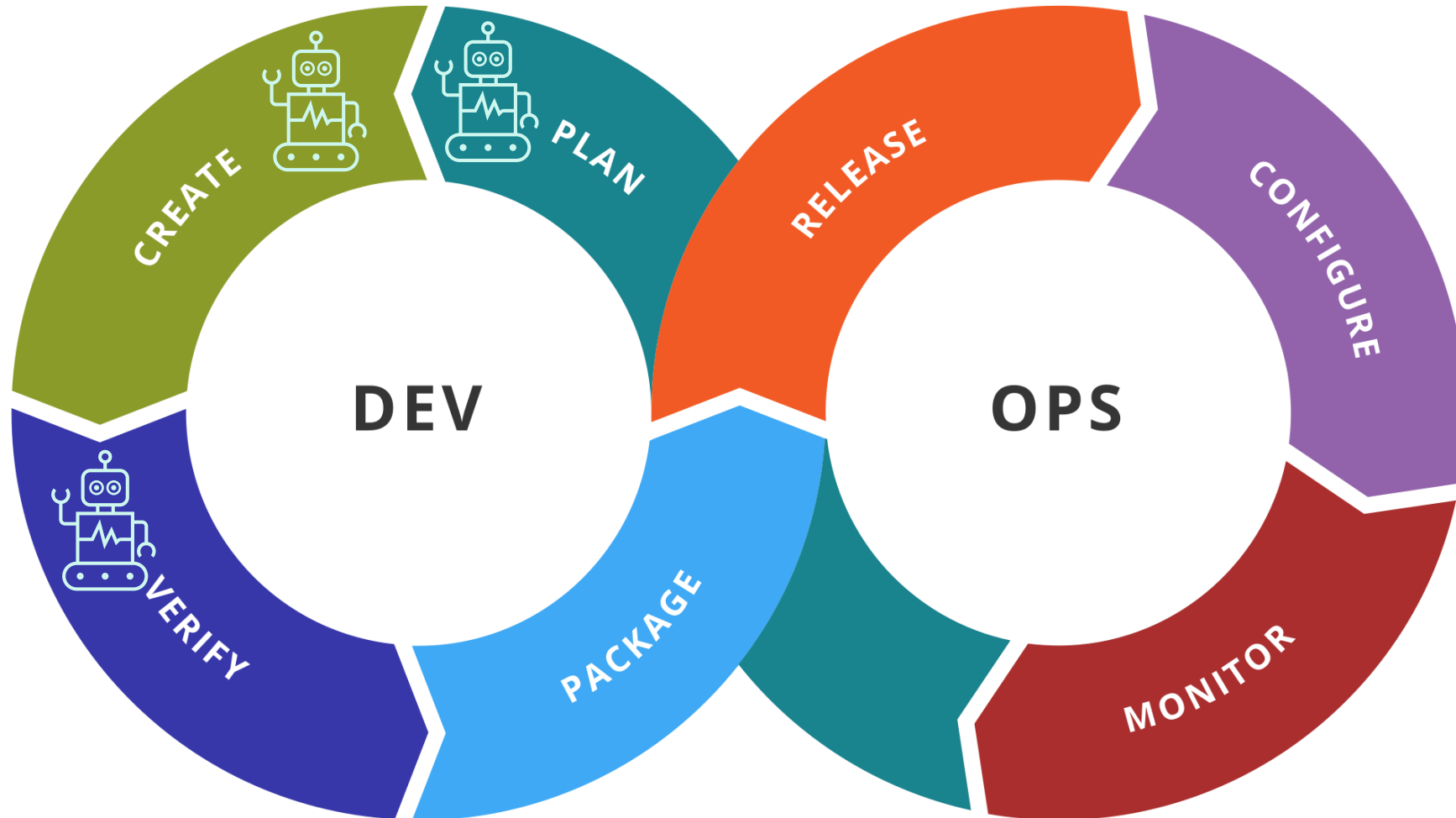# What's Next? – Shifting Left?

# What's Next if I am …

## A researcher:

AI will increasingly interact with people – **human subject studies** are key

- Enabling **productive collaboration** with AI agents

- Start with: how do developers currently collaborate & communicate?

- Then: how will that **change**, as more and more SE is done by AI?

- Managing **ownership and responsibilities** in AI-generated/maintained projects

# What's Next if I am ...

## A researcher:

To make AI effective, we need **new metrics and benchmarks**

- How to evaluate a comment? A design document? An entire PR?

- We need **next-generation benchmarks**, possibly LLM-powered.

- Building **frameworks of performance** to support developers, end-users

- As LLMs enter UX, we need design patterns for agents, test suites for UIs

# What's Next if I am …

## A researcher:

To democratize code and AI, we need to **rethink programming**

- Enabling code generation for **8 billion non-programmers**

- **Supporting learning**, debugging, maintenance for end-users

- Bringing together programmers and non-programmers
  - New development environments. Developers as a Service.

# What's Next if I am ...

## A CS student/professional:

Learn the **tools**

**Commercial**: Copilot, ChatGPT, Bard, Claude, ...

**Open-Source**: InCoder, StarCoder (bidirectional context), CodeGen (strong on Python), (Code)LLaMa (particularly large)

This list **changes constantly**, so important to stay up to date

Not a one-way street: **add value** with discernment, planning
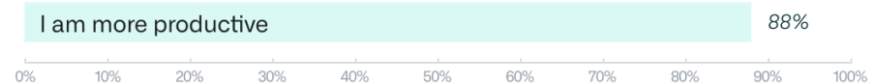
# What's Next if I am ...

## A CS student/professional:
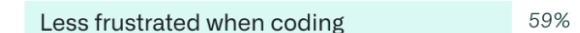
Shift the **emphasis** of your work

- Let AI do the **tedious** stuff

- Matplotlib's parameters have ridiculous names that you can't remember? Don't try, use AI

- Need a boilerplate website? Go ask a conversational LLM

- The set of boring things AI does easily is growing fast

**When using GitHub Copilot...**

**Perceived Productivity**

I am more productive — 88%

0%  10%  20%  30%  40%  50%  60%  70%  80%  90%  100%

**Satisfaction and Well-being**[*]

Less frustrated when coding — 59%

# What's Next if I am ...

## A CS student/professional:

Prepare for AI to do most of the coding

Prioritize **people & design** skills

- AI won't tell you what to build next, or who uses your products and why

Keep track of **what's important**

- AI is fueled by our past achievements. It changes what we think is hard & meaningful

- Keeping up with new tools helps you calibrate

---

**Geoffrey Hinton**
@geoffreyhinton

Caterpillars extract nutrients which are then converted into butterflies. People have extracted billions of nuggets of understanding and GPT-4 is humanity's butterfly.

4:27 PM · Mar 14, 2023 · **438.1K** Views

# What's Next if I am ...

**A non-programmer:**

A near-term future where almost anyone can **almost code**

- Some programming will be like using a screwdriver

- Other times it's more like fixing a car

- Not always obvious which it's going to be

# Beyond Code Generation

# Towards Next-Generation AI for SE

Vincent J. Hellendoorn

December 3rd, 2023

MAPS @ FSE

vhellendoorn@google.com