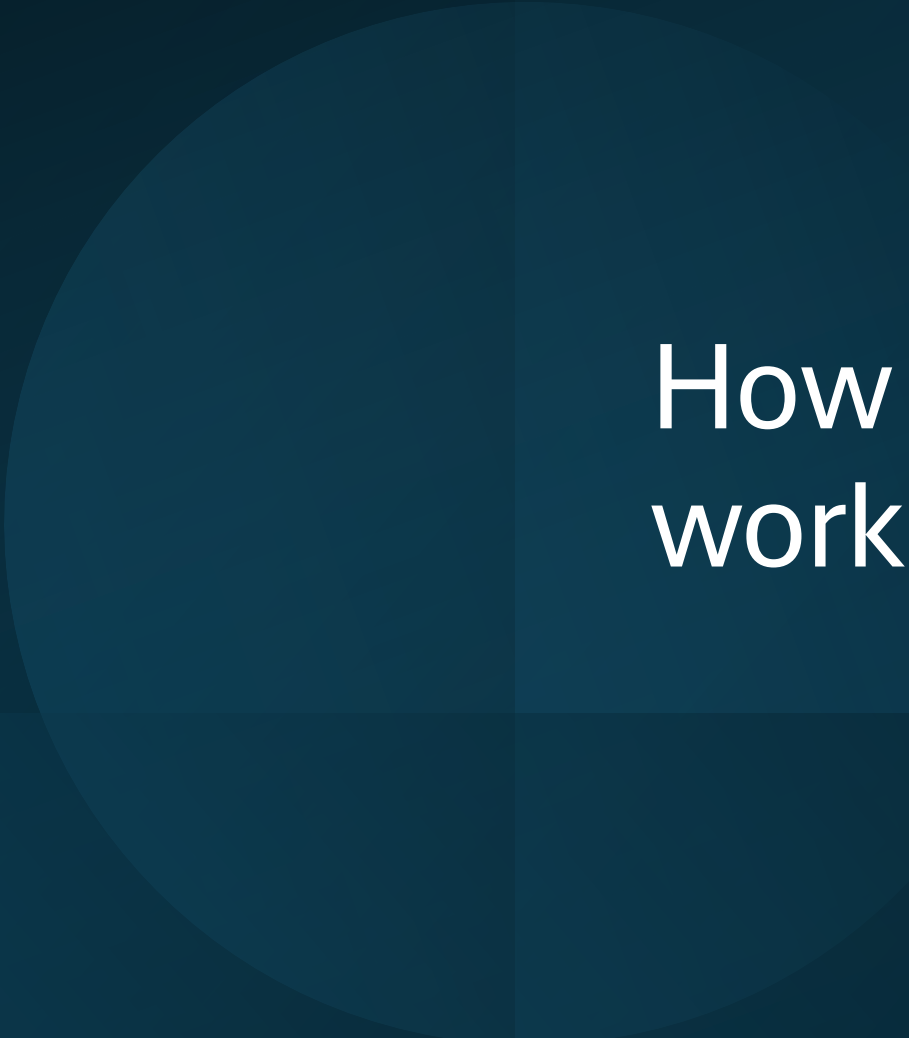# Headaches of shipping AI in products

Austin Z. Henley  •  Carnegie Mellon University  •  9/19/2024
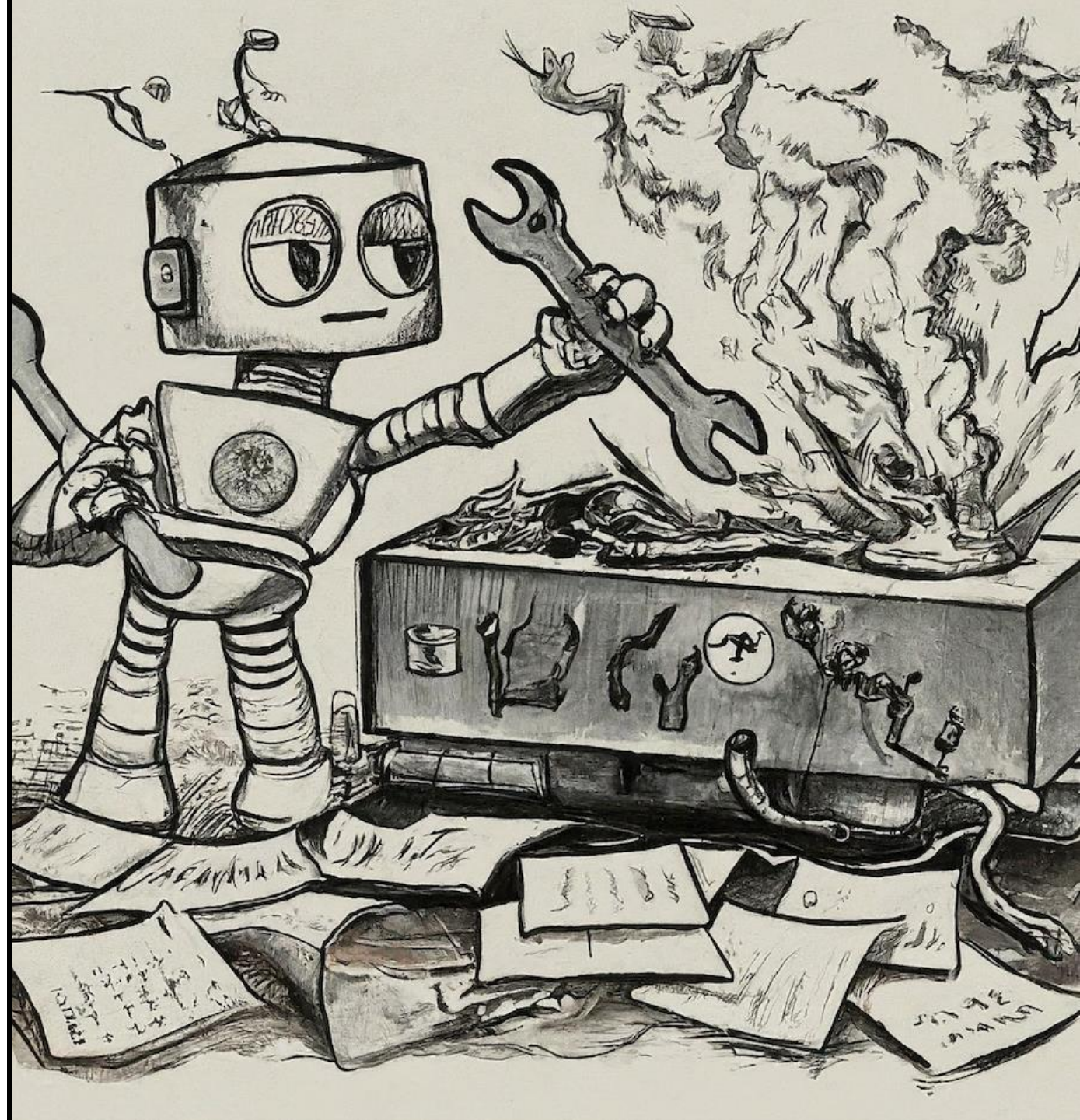
How many of you want to
work on AI?

What AI features do you find useful now?

# What is different for developing LLM-powered features?

- Non-deterministic
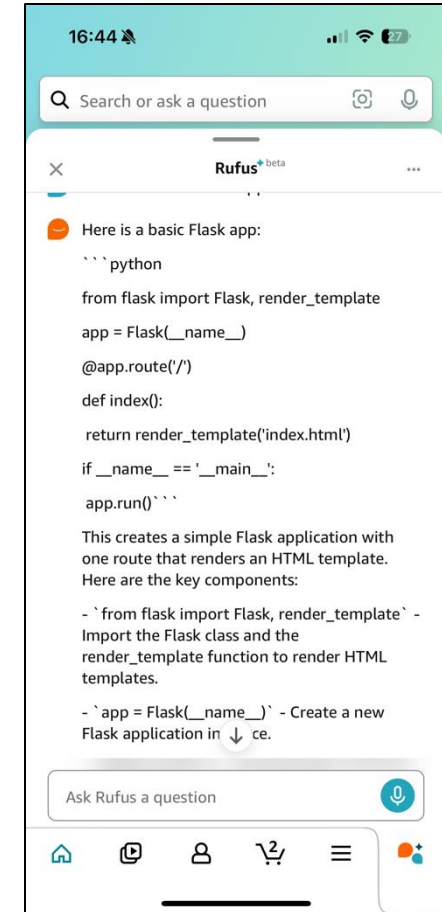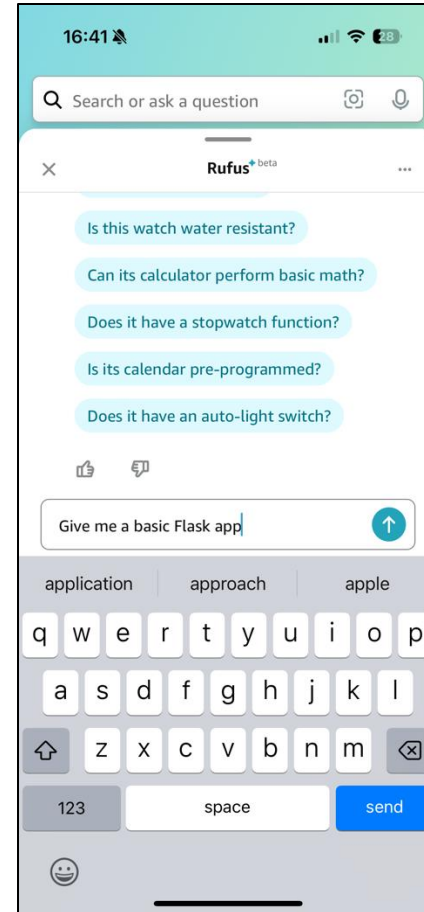- Costly
- Slow
- Safety
- Privacy
- Testing

- What is the user problem we are trying to solve?

- When is the AI good enough?

- How do you know when it goes off the rails?



*pRoMpT eNgInEeR*

# Ship it!

- When $1,000,000,000 is not enough
- We built this cool thing… can we ship it?
- Ok, we shipped… why aren't people using it?
- Ok, we shipped… how do we know if it is helpful?
- Let's try 50 different UIs
- How do we enable other developers to ship AI?

Before ChatGPT...

When $1,000,000,000 is not enough

# Our Mission

Providing insight to the instructors: classroom analytics. 👩‍🏫

AI Coding Tutor that is indistinguishable from human instructors. 👩‍🎓

Individual student's learning journey

Aggregated (class-wide) performance and pain-points

Efficiently give fair feedback

**Solution:** A dashboard that provides **quick summary** of students' performance and an interface to **provide mass feedback at scale**

# Goal: Assist instructors to get better insight of students' progress



While students are our audience, teachers are our **customers**

Our tools must be built to improve the learning and teaching experience.

# Classroom analytics: Dashboard for intro to Python course

# We built this cool thing… can we ship it?

C2

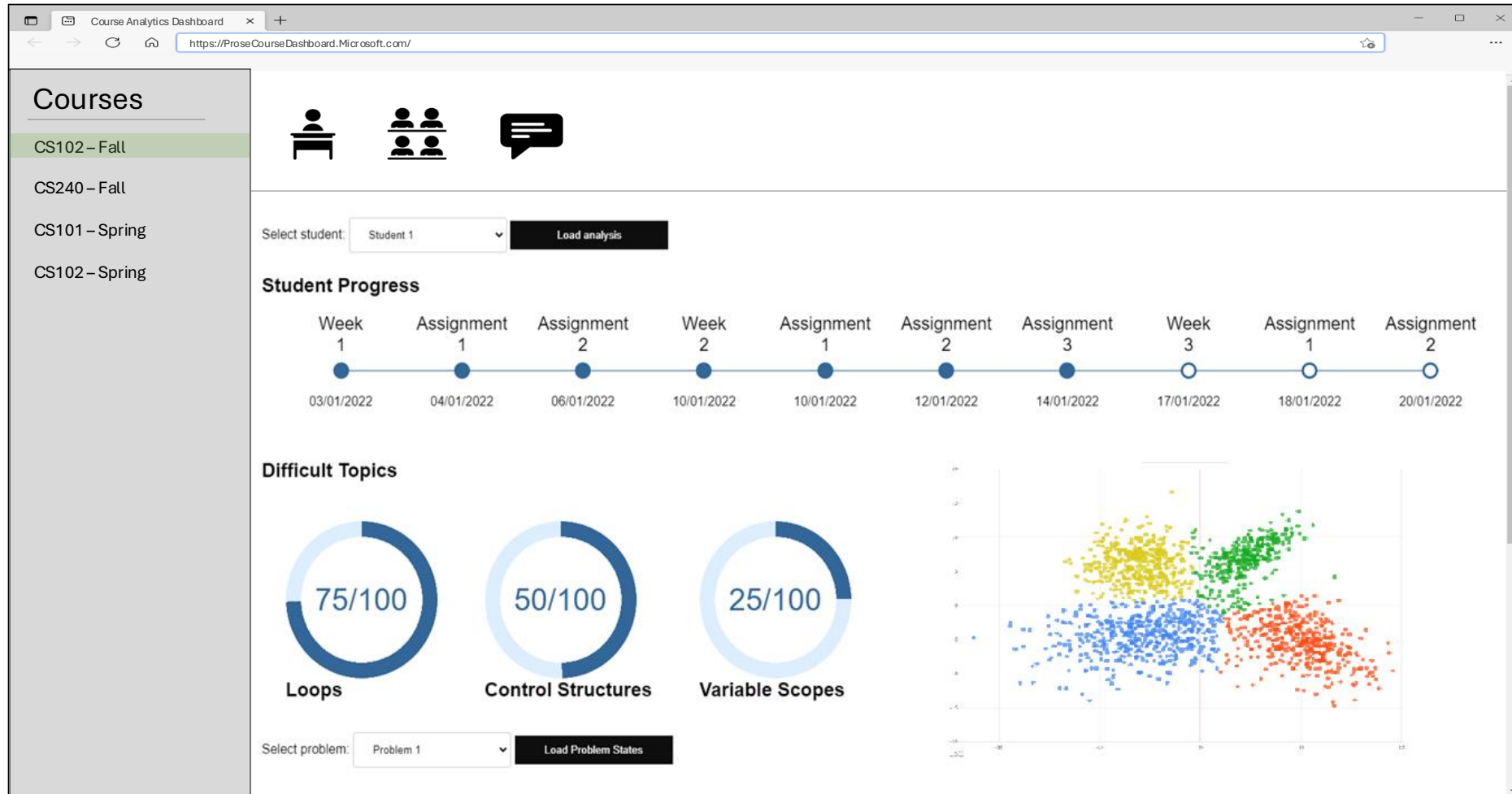| | A | B | C | D | E | F | G | H | I | J | K | L |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | First name | Last name | Initials | | | | | | | | | |
| 2 | Nathan | Cooper | | | | | | | | | | |
| 3 | Johnny | Phillips | | | | | | | | | | |
| 4 | Edward | Cox | | | | | | | | | | |
| 5 | Kyle | Howard | | | | | | | | | | |
| 6 | Christopher | Nguyen | | | | | | | | | | |
| 7 | Noah | Rivera | | | | | | | | | | |
| 8 | George | Patel | | | | | | | | | | |
| 9 | Denise | Roberts | | | | | | | | | | |
| 10 | Beverly | Moore | | | | | | | | | | |
| 11 | Kenneth | Bennet | | | | | | | | | | |
| 12 | Virginia | Sanchez | | | | | | | | | | |
| 13 | George | Lopez | | | | | | | | | | |
| 14 | Lisa | James | | | | | | | | | | |
| 15 | Jacqueline | Gonzales | | | | | | | | | | |
| 16 | Kyle | Myers | | | | | | | | | | |
| 17 | Andrew | Moore | | | | | | | | | | |
| 18 | Jason | Brooks | | | | | | | | | | |
| 19 | | | | | | | | | | | | |
| 20 | | | | | | | | | | | | |
| 21 | | | | | | | | | | | | |
| 22 | | | | | | | | | | | | |

Untitled-1.ipynb · df [DW] ✕

☐ Search

**OPERATIONS**

New column by example

Automatically create a column when a pattern is detected from the examples you provide. Powered by Microsoft Flash Fill.

Target columns

Name, Year_0

Derived column name

derivedCol

Apply   Discard

**DATA SUMMARY**

| | |
|---|---|
| Data shape | 16,598 rows x 13 columns |
| Columns | 13 |
| Rows | 16,598 |
| Missing values (by column) | 871 |

Export to notebook   Export as file   Copy all code   Report an issue

16598 rows x 13 columns   Go to column   Editing

| # index | # Rank | 🅰 Name | 🅰 Platform | 🅰 Year_0 | ⊞🅰 derivedCol ✓ | 🅰 Year_1 | 🅰 Genre |
|---|---|---|---|---|---|---|---|
| Missing: 0 (0%) | Missing: 0 (0%) | Missing: 0 (0%) | Missing: 271 (2%) | Missing: 271 (2%) | Missing: 271 (2%) | Missing: |
| Distinct: 16598 (100%) | Distinct: 16598 (100%) | Distinct: 11493 (69%) | Distinct: 31 (<1%) | Distinct: 39 (<1%) | Distinct: 822 (5%) | Distinct: 1 (<1%) | Distinct: 1 |
| | 11493 Distinct values | DS: 13% PS2: 13% PS3: 8% Other: 66% | 2009: 9% 2008: 9% 2010: 8% Other: 74% | 822 Distinct values | 0: 98% | Action: Sports: Misc: Other: |
| Min 1 · · · Max 16600 | | | | | | | |
| 0 | 1 | Wii Sports | Wii | 2006 | ✓ W_06 | 0 | Sports |
| 1 | 2 | Super Mario Bros. | NES | 1985 | S_85 | 0 | Platform |
| 2 | 3 | Mario Kart Wii | Wii | 2008 | M_08 | 0 | Racing |
| 3 | 4 | Wii Sports Resort | Wii | 2009 | W_09 | 0 | Sports |
| 4 | 5 | Pokemon Red/Pokemon Blue | GB | 1996 | P_96 | 0 | Role-Playing |
| 5 | 6 | Tetris | GB | 1989 | T_89 | 0 | Puzzle |
| 6 | 7 | New Super Mario Bros. | DS | 2006 | N_06 | 0 | Platform |
| 7 | 8 | Wii Play | Wii | 2006 | W_06 | 0 | Misc |
| 8 | 9 | New Super Mario Bros. Wii | Wii | 2009 | N_09 | 0 | Platform |
| 9 | 10 | Duck Hunt | NES | 1984 | D_84 | 0 | Shooter |
| 10 | 11 | Nintendogs | DS | 2005 | N_05 | 0 | Simulation |
| 11 | 12 | Mario Kart DS | DS | 2005 | M_05 | 0 | Racing |
| 12 | 13 | Pokemon Gold/Pokemon Silver | GB | 1999 | P_99 | 0 | Role-Playing |
| 13 | 14 | Wii Fit | Wii | 2007 | W_07 | 0 | Sports |
| 14 | 15 | Wii Fit Plus | Wii | 2009 | W_09 | 0 | Sports |
| 15 | 16 | Kinect Adventures! | X360 | 2010 | K_10 | 0 | Misc |
| 16 | 17 | Grand Theft Auto V | PS3 | 2013 | G_13 | 0 | Action |

**CLEANING STEPS**

1 Load data from variable

2 Change column type   'Year'

3 Split text   'Year'

4 **New column by example**   2 columns

☐ Preview code for all steps

PROBLEMS   OUTPUT   DEBUG CONSOLE   TERMINAL   PORTS   JUPYTER   **DATA WRANGLER**

4 *New column by example*

```
1   # Derive column 'derivedCol' from columns: 'Year_0', 'Name'
2   # Transform based on the following examples:
3   #    Name           Year_0    Output
4   # 1: "Wii Sports"   "2006" => "W_06"
5   df.insert(4, "derivedCol", df["Name"].str[:1] + "_" + df["Year_0"].str[2:])
```

⚡ Previewing

Apply   Discard

⊗ 0 ⚠ 0   ⊘ 0   ✎ Data Wrangler: Editing   ⊞ Operation preview: Derive column 'derivedCol' from columns: 'Year_0', 'Name'   Pandas 2.2.2   idle

Ok, we shipped…
why aren't people
using it?

Ok, we shipped... how do we know if it is helpful?

How do we know when Excel Copilot is helping users? Not helping?

# Let's try 50 different UIs

```
return ((fMin - 32) * (5.0 / 9.0));
```

**edit 1**

```
return (FtoC(fMin));
```

```
return (fMax - 32) * (5.0 / 9.0);
```

**edit 2**

```
return FtoC(fMax);
```

```
return (fAve - 32) * (5.0 / 9.0);
```

```
141         {
142             fTot += f;
143         }
144
145         fAve = fTot / fTemps.Count;
146
147         return (fAve - 32) * (5.0 / 9.0);
148     }
149

    0 references | 0 changes | 0 authors, 0 changes
150     static double SumTempsInC(List<double> fTemps)
151     {
152         double fTotal = 0;
153         foreach (double f in fTemps)
154         {
155             fTotal += f;
```

```
return (fAve - 32) * (5.0 / 9.0);
```

💡 ▾

★ IntelliCode suggestion based on recent edits: FtoC(fAve) ▸   Apply suggestion

Ignore suggestions like this

reference | Gustavo Soares, 288 days ago | 1 author, 1 change

```
tatic double AveCTempInF(List<double> cTemps)


    double cTot = 0;

    double cAve:
```

```
...
return (fAve - 32) * (5.0 / 9.0);
return FtoC(fAve);
}
...
```

```csharp
class Obstacle
{
    public ObstacleType type { get; set; }
    public int XPos { get; set; }
    public int YPos { get; set; }
    public double XVelocity { get; set; }
    public double YVelocity { get; set; }
    private bool IsValid { get; set; }

    public Obstacle (ObstacleType type, int xPos, int yPos, double xVelocity, double yVelocity, bool isValid)
    {
        Type = type;
        XPos = xPos;
        XVelocity = xVelocity;
        YVelocity = yVelocity;
        IsValid = isValid;
    }
}
```

Tab accept ...

(a) `if (obstacle != null && obstacle.IsValid)` ⟶ `ObjectNotNullAndValid(obstacle)` `Tab accept | ...`

(b) `if (ObjectNotNullAndValid(obstacle) != null && obstacle.IsValid)` `Tab accept | ...`

(c) `if (ObjectNotNullAndValid(obstacle))` `Tab accept | ...`

(d) `if (obstacle != null && obstacle.IsValid)` ⟶ `ObjectNotNullAndValid(obstacle)` `Tab accept | ...`

## (a)

```
obstacleObjects.ForEach(obstac.  Tab  to accept
{
    if (obstacle != null && obstacle.IsValid)
    {
        if (obstacle.XPos < 0 || obstacle.YPos < 0 || obstacle.XPos > 200 || obstacle.YPos > 200)
        {
            obstacle.IsValid = false;
        }
    }
    if (ObjectNotNullAndValid(obstacle))
    {
        if (obstacle.IsOutOfFrame())
        {
            obstacle.SetInvalid();
        }
    }
});
```

(a)

## (b)

```
obstacleObjects.ForEach(obstacle =>
{
    if (obstacle != null && obstacle.IsValid)
    {
        if (obstacle.XPos < 0 || obstacle.YPos < 0 || obstacle.XPos > 200 || obstacle.YPos > 200)
        {
            obstacle.IsValid = false;
        }
    }
```

Refactor Preview.  Tab  to accept   Ctrl + ]  for next suggestion   Ctrl + [  for previous suggestion   ✕

```
    if (ObjectNotNullAndValid(obstacle))
    {
        if (obstacle.IsOutOfFrame())
        {
            obstacle.SetInvalid();
        }
    }
```

BirdObstacle

(b)

## (c)

```
obstacleObjects.ForEach(obstacle =>
{
    if (obstacle != null && obstacle.IsValid)
    {
        if (obstacle.XPos < 0 || obstacle.YPos < 0 || obstacle.XPos > 200 || obstacle.YPos > 200)
        {
            obstacle.IsValid = false;
        }
    }
```

Refactor Preview.  Tab  to accept   Ctrl + ]  for next suggestion   Ctrl + [  for previous suggestion   ✕

```
    if (ObjectNotNullAndValid(obstacle))
    {
        if (obstacle.IsOutOfFrame())
        {
            obstacle.SetInvalid();
        }
    }
```
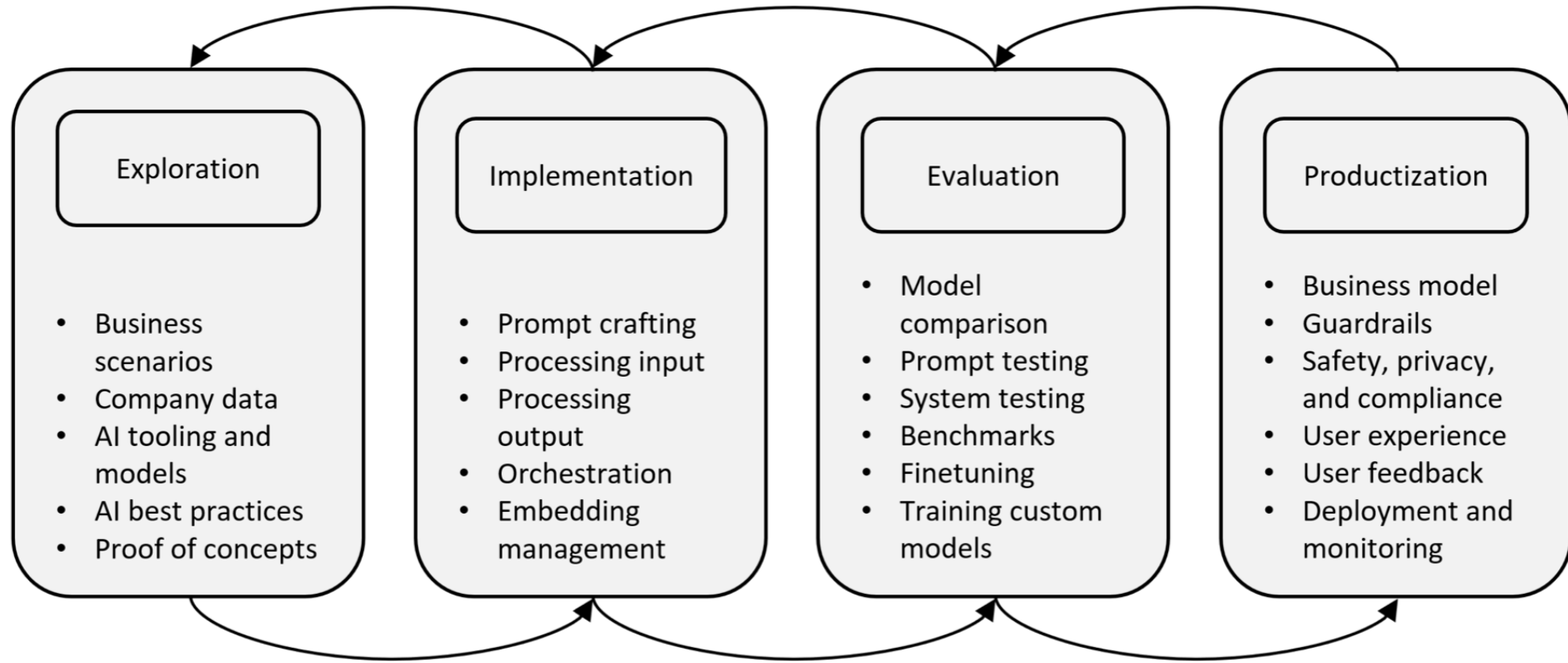
BirdObstacle

(c)

## (d)

```
obstacleObjects.ForEach(obstacle =>
{                    Ctrl + .   see suggestion
    if (obstacle != null && obstacle.IsValid)
    {
        if (obstacle.XPos < 0 || obstacle.YPos < 0 || obstacle.XPos > 200 || obstacle.YPos > 200)
        {
            obstacle.IsValid = false;
        }
    }
});
```

(d)

# How do we enable other developers to ship AI?

People were using unit tests to…

# Lessons, maybe



Are you a solution looking for a problem?



Big companies have many constraints—many have nothing to do with how good your tech is
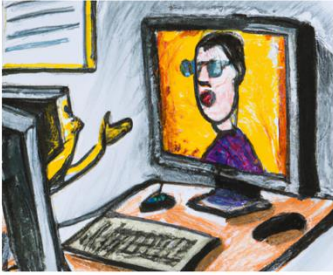


Innovator's dilemma?

# Window 1

**Austin Z. Henley**
Associate Teaching Professor
Carnegie Mellon University

azhenley@cmu.edu
@austinzhenley
github/AZHenley

## Natural language is the lazy user interface

1/27/2023



See the discussion of this post on *Hacker News*.

ChatGPT has kicked off a frenzy. It is all anyone in the tech world is talking about it seems. Startups are popping up left and right. Big companies are rapidly releasing ChatGPT-like features integrated in their products.

People are anticipating that large language models are going to revolutionize the world.

And maybe they will.

But a chat bot won't.

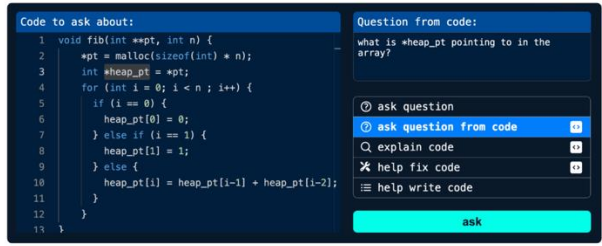Expecting users to primarily interact with software in natural language is *lazy*.

---

# Window 2

**Austin Z. Henley**
Associate Teaching Professor
Carnegie Mellon University

azhenley@cmu.edu
@austinzhenley
github/AZHenley

## CodeAid: A classroom deployment of an LLM-based programming assistant

5/19/2024



This post was co-written with *Majeed Kazemitabaar*, who led this project. Majeed is a PhD student in CS at the University of Toronto who has been researching the educational impact and utility of LLMs in computing education. We summarize our recent CHI'24 paper, *"CodeAid: Evaluating a Classroom Deployment of an LLM-based Programming Assistant that Balances Student and Educator Needs"*. See the *paper* for more details.

See the discussion of this post on *Hacker News*.

---

LLM-powered tools like ChatGPT can assist students that need help in programming classes by explaining code and coding concepts, generating fixed versions of incorrect code, providing examples, suggesting areas of improvement, and even writing entire

---

# Window 3

**Austin Z. Henley**
Associate Teaching Professor
Carnegie Mellon University

azhenley@cmu.edu
@austinzhenley
github/AZHenley

## Exploring 50 user interfaces for AI code suggestions

5/7/2024

*This post is a summary of our ICSE-SEIP'23 paper, "Towards More Effective AI-Assisted Programming: A Systematic Design Exploration to Improve Visual Studio IntelliCode's User Experience". See the paper for more details. Thanks to Priyan Vaithilingam for leading this project!*

AI code suggestions in code editors, such as Copilot and Visual Studio IntelliCode, are fairly common place now.



What is the optimal way to present these code suggestions to users? Showing the suggestion as inline gray text is common. But are there better ways? What if it is a multi-line change that modifies existing code?

These are the questions we wanted to answer.

### Methodology

We iteratively explored 50 or so designs for inline code change interfaces in Visual Studio. We filtered our designs down to 19 which we then tested in a series of 7 lab