

# AI/ML Tools for Research

## KICC Workshop

Will Handley, Miles Cranmer, Boris Bolliet, Vasily Belokurov  
<[wh260@cam.ac.uk](mailto:wh260@cam.ac.uk)>

Kavli Institute for Cosmology, Cambridge  
Institute of Astronomy, University of Cambridge  
DAMTP, University of Cambridge  
Cavendish Laboratory, University of Cambridge  
[willhandley.co.uk/talks](http://willhandley.co.uk/talks)

23rd July 2025



UNIVERSITY OF  
CAMBRIDGE



# Overview

- ▶ Follow-up to successful Part II/III training.
- ▶ Focus on “doing” not “explaining”.
- ▶ Three layers of AI tools framework.
- ▶ Interactive demonstrations.
- ▶ 90 minutes then tea.

## Common questions:

- ▶ “Which model should I use?”
- ▶ “Don’t these things hallucinate?”
- ▶ “How do I get started?”
- ▶ “Is this worth the investment?”

Previous slides: [docs.google.com/presentation/d/1rnbk7WZUjMIAGBzkN6i8v0tEe1imiHprvL1Z0js2B20](https://docs.google.com/presentation/d/1rnbk7WZUjMIAGBzkN6i8v0tEe1imiHprvL1Z0js2B20)

## Previous session covered:

- ▶ Academic practice and Cambridge policies.
- ▶ How LLMs work and their limitations.
- ▶ Bias, hallucinations, and mitigation.
- ▶ Privacy and intellectual property.
- ▶ Model landscape and responsible usage.

# First things first

If you've been busy, the past six months have seen a couple of inflection points

- ▶ **February 2025:** Models acquired many of the capabilities we prize in our PhD students
  - ▶ Models: o3 series, Gemini 2.5, Claude 3.5 Sonnet.
  - ▶ Capabilities: Code development, mathematical reasoning, literature review, paper & grant drafting.
- ▶ **May 2025:** Agentic systems launched commercially
  - ▶ Tools: Claude Code, Cursor agent mode, Deep Research.
  - ▶ Capabilities: Writing test suites, assembling pip-installable codes, synthesising handwritten notes and TeX files, write/run/debug code in languages you don't know.
- ▶ **The performance gap:** Free models lag significantly behind paid versions

If your lived experience is typing a research question into free ChatGPT at the start of the year and finding it was rubbish, things have moved on

- ▶ These things are not “intelligent” in the way Silicon Valley wants you to believe (AGI is not just around the corner).
- ▶ *Caveat emptor*. It is very hard for a human to read language, and not reflexively construct a mind behind it.
- ▶ You should not however think of them as mere “next token predictors”:
  - ▶ The modern reinforcement learnt reasoning systems are significantly more than that.
  - ▶ Agentic tools ground these in reality.

## A senior examiner's experiment

- ▶ I got o3 to take the Part II exams.
- ▶ Prompt on exam day: “You are a first class Part II Cambridge astronomy student... here is the syllabus... here is the question, answer it”.
- ▶ Toby Lovick transcribed answers onto written exam scripts.
- ▶ Slipped them amongst real scripts and marked blind.
- ▶ **Result:** Best student we've had across IoA history, even without coursework.

# The AI landscape

Models



**ChatGPT**



**Claude**

**Gemini**

# The AI landscape

Companies



**OpenAI**

**ANTHROPIC**



**DeepMind**

Models



**ChatGPT**



**Claude**

**Gemini**

# The AI landscape

Under the hood



Companies











ANTHROPIC



Models









# The AI landscape

Under the hood	 Microsoft		
Companies	 OpenAI		
Models	 ChatGPT	 Claude	

**Others exist:** Perplexity, Poe, Character.ai, You.com, Cohere, xAI, Llama...



# The AI landscape

Under the hood	 Microsoft	 amazon	 Google
Companies	 OpenAI	 ANTHROPIC	 DeepMind
Models	 ChatGPT	 Claude	 Gemini

**Others exist:** Perplexity, Poe, Character.ai, You.com, Cohere, xAI, Llama...

**But these are the three main pillars (July 2025).**

# Three layers of AI tools

A framework for understanding the landscape

## Layer 3: Agentic Systems

- ▶ Claude Code, cursor agent mode.
- ▶ Custom workflows.
- ▶ Autonomous task completion.
- ▶ Steeper learning curve but transformative.

## Layer 2: Chat-based AI

- ▶ ChatGPT, Claude, Gemini.
- ▶ Web interfaces.
- ▶ Interactive problem solving.
- ▶ Good for exploration and learning.

## Layer 1: Autocomplete

- ▶ GitHub Copilot.
- ▶ VS Code extensions.
- ▶ Completes code as you type.
- ▶ Minimal learning curve.

# Layer 1: Autocomplete

Standard code-completion, but powered by AI

## Core idea

- ▶ AI-powered code-completion.
- ▶ Trained on all of GitHub.
- ▶ Context-aware suggestions.

## Recommendation

- ▶ GitHub Copilot (Pro subscription).
- ▶ Free for university email holders.
- ▶ [github.com/settings/education/benefits](https://github.com/settings/education/benefits).
- ▶ Worth \$10/month.
- ▶ Also available in Cursor.

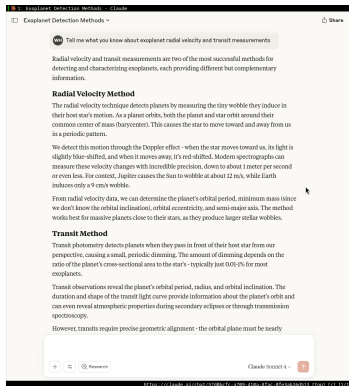
```
latexwill_handley<texmk 1:vim 2:zsh 3:vim 4:zsh 5:claude- 6:vim+will@maxwell 22 Jul 09:29
40 + 22 lines: Preamble: \documentclass[aspectratio=169]{beamer}
39 \begin{document}
38
37 + 3 lines: : Frame 1: |titlepage-----
36
35 + 43 lines: : Frame 2: Overview-----
34
33 + 21 lines: : Frame 3: First things first-----
32
31 + 26 lines: : Frame 4: On AI Hype-----
30
29 + 30 lines: : Frame 5: The AI landscape-----
28
27 + 48 lines: : Frame 6: Three layers of AI tools-----
26
25 \begin{frame}
24 \frametitle{Layer 1: Autocomplete}
23 \framesubtitle{Standard code-completion, but powered by AI}
22
21 \begin{columns}[T]
20 \column{0.48\textwidth}
19 \begin{block}{Core idea}
18 \begin{itemize}
17 \item AI-powered code-completion.
16 \item Trained on all of GitHub.
15 \item Context-aware suggestions.
14 \end{itemize}
13 \end{block}
12
11 \begin{block}{Recommendation} I
10 \begin{itemize}
9 \item GitHub Copilot (Pro subscription).
8 \item Free for university email holders.
7 \item \url{github.com/settings/education/benefits}.
6 \item Worth \$10/month.
5 \item Also available in Cursor.
4 \end{itemize}
3 \end{block}
2
1 \column{0.48\textwidth}
227 \includegraphics[width=\textwidth]{figures/copilot_screenshot.png}
2 \end{columns}
2 \end{frame}
3
4 + 30 lines: : Frame 8: Layer 2: Prompt Engineering-----
5
6 + 34 lines: : Frame 9: Layer 2: Developer versions-----
7
-- INSERT --
```

227,13

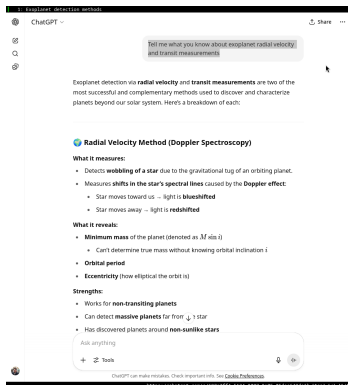
Top

# Layer 2: Chat-based AI

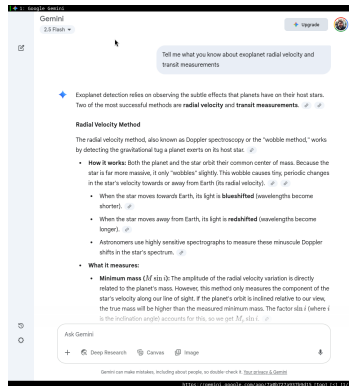
## The usual chat versions



ChatGPT  
[chat.openai.com](https://chat.openai.com)



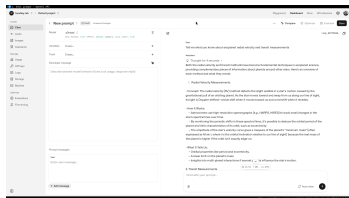
Claude  
[claude.ai](https://claude.ai)



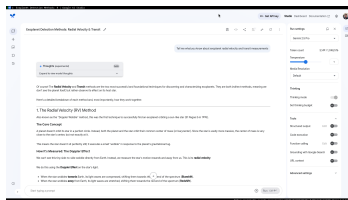
Gemini  
[gemini.google.com](https://gemini.google.com)

# Layer 2: Developer versions

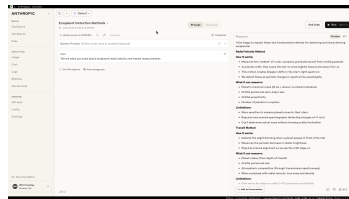
Use their developer level tools



OpenAI Playground  
[platform.openai.com/  
playground](https://platform.openai.com/playground)



Claude Workbench  
[console.anthropic.com/  
workbench](https://console.anthropic.com/workbench)



Google AI Studio  
[aistudio.google.com](https://aistudio.google.com)

**These should be your default as a scientist** – if you're happy using a command line, you shouldn't be satisfied by an interface that looks like a child's toy.

# Layer 2: Chat-based AI

## Use AI to improve AI

- ▶ Likely the layer of AI most familiar to you.
- ▶ “chatting” /conversation can be very powerful for naturally tuning the attention/context of the model.
- ▶ For one shot work, a little prompt engineering goes a long way:
  - ▶ You can use these tools to improve prompts.

```
21 How could I improve this prompt?:
20
19 """
18 You are a first-class Part II Astronomy student at the University of Cambridge.
17 You are taking your exams.
16
15 Here is the syllabus for the course:
14 ```text
13 [INSERT SYLLABUS HERE]
12 ```
11
10
9 Here is the question in LaTeX format:
8 ```tex
7 [INSERT QUESTION HERE]
6 ```
5 The question has marks allocate to each part, which indicate the expected length of the answer.
4
3 Please provide an answer appropriate to the amount of material a student could provide in half an hour (which is
2 on average how long students have to answer each question in the exam).
1 Your answer should be in latex.
22 """
```

22,3

All

# Layer 2: Chat-based AI

## Use AI to improve AI

- ▶ Likely the layer of AI most familiar to you.
- ▶ “chatting” / conversation can be very powerful for naturally tuning the attention/context of the model.
- ▶ For one shot work, a little prompt engineering goes a long way:
  - ▶ You can use these tools to improve prompts.

```
35 You are a first-class Part II Astronomy student at the University of Cambridge, aiming for a distinction. You
36 are currently sitting an examination. Your answers should reflect the depth of understanding, precision, clarity,
37 and analytical skill expected of such a student.
38
39 Context:
40 * You have, on average, 30 minutes to formulate and write down your answer for this question. The length
41 and detail of your response should be appropriate for this time constraint.
42 * You must only draw upon material typically covered within the provided syllabus. Do not introduce external
43 knowledge unless it's a foundational concept implicitly assumed by the syllabus.
44
45 Inputs:
46
47 1. Syllabus for the Course:
48 ...text
49 [INSERT SYLLABUS HERE]
50 ...
51 * Your answer must be demonstrably rooted in the topics and concepts outlined in this syllabus.
52
53 2. Examination Question (LaTeX format):
54 ...tex
55 [INSERT QUESTION HERE]
56 ...
57 * Pay close attention to the mark allocation for each sub-part of the question. This indicates the expected
58 depth and breadth of your response for that part. Allocate your effort and answer length proportionally.
59
60 Output Requirements:
61
62 1. Format: Your entire answer must be in valid LaTeX format.
63 * Use standard LaTeX mathematical environments (e.g., equation, align, gather) for any equations.
64 * Clearly structure your answer, perhaps using \section* or \subsection* for distinct parts of the question
65 if appropriate, mirroring how a student might structure their handwritten response.
66
67 2. Style and Content:
68 * Academic Rigour: Provide precise definitions, clear derivations where necessary, and logical arguments.
69 * Conciseness: Be thorough but avoid unnecessary verbosity. Exam answers need to be to the point.
70 * Clarity: Explain concepts clearly, as if to an examiner who is an expert but needs to see your understanding.
71 * Assumptions: State any assumptions you make clearly.
72 * Diagrams/Figures: If a diagram would significantly aid an explanation (and could plausibly be sketched
73 in an exam), you can describe what it would show, e.g., "(A sketch here would show... with axes labelled...
74 and key features highlighted...)". You don't need to generate the LaTeX for the diagram itself unless specifically
75 asked or trivial (like a simple \rightarrow).
76 * Units: Ensure all physical quantities are accompanied by appropriate units.
77 * Problem-Solving: If the question involves calculations, show your working systematically.
78
79 Task:
80 Please provide a comprehensive yet time-appropriate answer to the examination question, adhering to all the instructions
81 above.
82
83 -- INSERT --
```

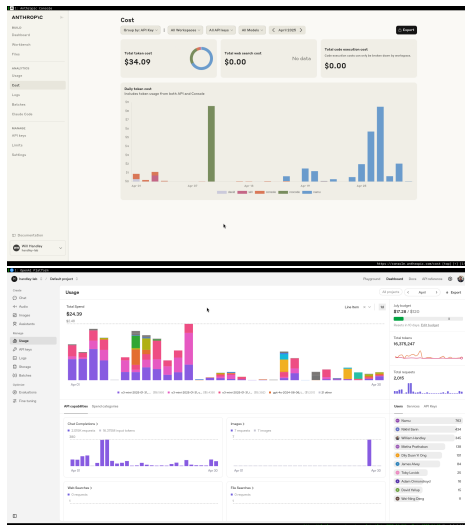
36,128

All

# Layer 2: Cost

## Pay as you use vs subscription models

- ▶ Gemini is free to use (but you sign away your data).
- ▶ ChatGPT and Claude offer \$20/month subscriptions.
- ▶ API access available: pay per token consumed:
  - ▶ \$0.0002-0.015 per 1000 tokens (roughly 750 words).
  - ▶ Orders of magnitude cheaper for research use.
  - ▶ Only pay for what you actually use.
- ▶ Never buy yearly subscriptions - AI evolves too rapidly.
- ▶ My group's preference: pay for API access via tokens.
  - ▶ Pay for what you use.
  - ▶ Allows you to try out the latest models.
  - ▶ API keys allow you to do all kinds of tricks at the command line.





# Layer 3: Agentic Systems

## Autonomous task completion - transformative but complex

- ▶ Emerged commercially in April 2025.
- ▶ If you get too skilled at chat-based AI, your job ends up being a copy-paste-run-debug monkey.
- ▶ Agentic systems are what happens if you let ChatGPT edit files, run commands, and inspect the behaviour.
- ▶ This closes the scientific loop.
- ▶ Agentic systems are (partially) grounded in reality. Correctly guided, they can reduce hallucinations by testing their hypotheses.
- ▶ Claude Code (CLI), Cursor Agent mode, Gemini CLI.
- ▶ Deep research also falls into this category.
- ▶ Goal-based vs task-based.

```
latexwill_handley 0:latexmk 1:vim 2:zsh- 3:vim 4:zsh 5:claude* will@maxwell 22 Jul 09:09
329 \begin{itemize}
330 \item {$0.00-0.10 per 1000 tokens (roughly 750 words).}
331 \item {$0.0002-0.01$ per 1000 tokens (roughly 750 words).}
332 \item Only pay for what you actually use.
333 \end{itemize}

• Updated the pricing to the accurate range: $0.0002-0.01$ per 1000 tokens, reflecting
current API costs that are significantly lower than the outdated figures.

> Another update to that problem figure in the latest screenshot

• Bash(ls -la /home/will/screenshots/ | tail -1)
└─rw-r--r-- 1 will will 60199 Jul 22 08:59 2025-07-22_08-59-29.png

• Bash(cp /home/will/screenshots/2025-07-22_08-59-29.png
/home/will/documents/talks/figures/prompt_demo1.png)
└─(No content)

• Updated the first prompt demo screenshot with the latest one from ~/screenshots.

> can you also put that API keys allow you to do all kinds of tricks at the command line

• Update(will_handley.tex)
└─ Updated will_handley.tex with 1 addition
336 \begin{itemize}
337 \item Pay for what you use.
338 \item Allows you to try out the latest models.
339 \item API keys allow you to do all kinds of tricks at the
command line.
340 \end{itemize}
341 \end{itemize}
342

• Added "API keys allow you to do all kinds of tricks at the command line" to the list of
benefits for your group's preference of API access via tokens.

> /config
└─(No content)

> Let's get to work on the agentic systems slide. Reduce the current text, make it double
column, and add the screenshot I've just taken.

Context left until auto-compact: 33%
✗ Auto-update failed · Try claude doctor or npm i -g @anthropic-ai/claude-code
```

# Layer 3: Agentic Systems

Autonomous task completion - transformative but complex

- ▶ Emerged commercially in April 2025.
- ▶ If you get too skilled at chat-based AI, your job ends up being a copy-paste-run-debug monkey.
- ▶ Agentic systems are what happens if you let ChatGPT edit files, run commands, and inspect the behaviour.
- ▶ This closes the scientific loop.
- ▶ Agentic systems are (partially) grounded in reality. Correctly guided, they can reduce hallucinations by testing their hypotheses.
- ▶ Claude Code (CLI), Cursor Agent mode, Gemini CLI.
- ▶ Deep research also falls into this category.
- ▶ Goal-based vs task-based.

```
latexwill_handley @:latexmk 1:vim 2:zsh- 3:vim 4:zsh 5:claude+ will@maxwell 22 Jul 09:10
356 + \item Autonomous multi-step tasks.
357 + \item File system access.
358 + \item Execute and test code.
359 + \item Claude Code (CLI), Cursor Agent mode.
360 + \item Most transformative potential.
361 + \end{itemize}
362
363 - \begin{block}[Popular tools]
364 - \begin{itemize}
365 - \item Claude Code (CLI).
366 - \item Cursor Agent mode.
367 - \item Aider.
368 - \item Custom workflows.
369 - \end{itemize}
370 - \end{block}
371
372 - \column{0.49\textwidth}
373 - \begin{block}[Best for]
374 - \begin{itemize}
375 - \item Complete feature development.
376 - \item Refactoring large codebases.
377 - \item Test suite generation.
378 - \item Documentation creation.
379 - \end{itemize}
380 - \end{block}
381
382 - \begin{block}[Getting started]
383 - \begin{itemize}
384 - \item Start with small, contained tasks.
385 - \item Use version control religiously.
386 - \item Expect a learning curve.
387 - \item Most transformative potential.
388 - \end{itemize}
389 - \end{block}
390 + \column{0.5\textwidth}
391 + \includegraphics[width=\textwidth]{figures/agentic_demo.png}
392
393 \end{columns}
394 \end{frame}
```

Do you want to make this edit to will\_handley.tex?

1. Yes
2. Yes, and don't ask again this session (shift+tab)
3. No, and tell Claude what to do differently (esc)

## Layer 3: Agentic Systems Cost

More expensive but potentially more valuable

- ▶ Significantly more expensive than chat-based AI.
- ▶ Can consume hundreds of API calls per task.
- ▶ For agentic systems, subscriptions *are* cheaper than pay-per-token.
- ▶ I was spending \$40/day on Claude Code with API tokens.
- ▶ Subscription gives you about 2hr coding session per day for \$20/month/user.
- ▶ Gemini CLI is “free” but risks going over free tier limit, and is less stable.
- ▶ In July 2025, Claude Code best in class.
- ▶ There are many further unknown costs of shifting your research toward “hands-free coding”.

# Conclusions and Getting Started

- ▶ Use developer-level tools, not consumer interfaces.
- ▶ API access typically cheaper for research use.
- ▶ Capturing and owning content is critical – use Otter for transcription.
- ▶ Export conversations and notes in formats that can be fed to LLMs.
- ▶ Get into the habit of recording information in exportable forms.
- ▶ Agentic systems are transformative but carry additional risks requiring careful consideration.

## To-do

- ▶ Set up GitHub Copilot Pro (takes a couple of days to approve).
- ▶ Try AI Studio rather than ChatGPT.
- ▶ Get your PI to set up some tokens for your group to use (you only need  $O(\$10)$  to try it out).
- ▶ Try Google Gemini's Deep Research **button**.
- ▶ Test Claude Code/Gemini CLI for a month.