# > How to code yourself into a mess

## (& hopefully out again)

@ Generative AI for Business Researchers
Tuesday 19th Nov 2024

**Dr. Gordon McDonald**
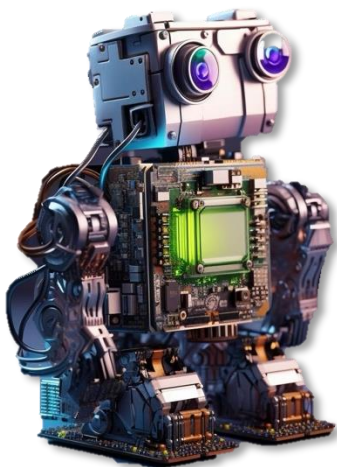Informatics Team Lead,
Sydney Informatics Hub

THE UNIVERSITY OF
SYDNEY
—
**Sydney
Informatics Hub**

informatics.sydney.edu.au/

# You're trusting this guy

- He's doing his best
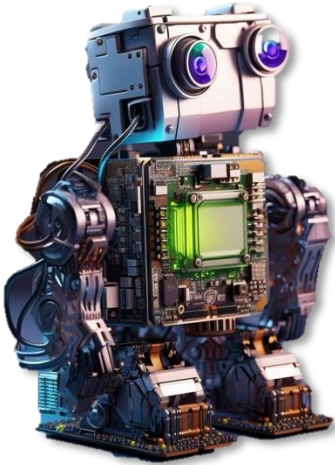
- But mistakes happen.

- Maybe you should pay him more

# I want some instant data science

Use an easy interface like ChatGPT, Anthropic…!
(USyd-approved=Copilot)

1. you'll have to pay ~US$20/mo
2. your data and code is going god-knows-where…
    - make sure it's only public data and not private info!
3. ask, be very specific
4. It has limited access to the internet, limited compute, limited file size upload.

# sudo make me a sandwich

# I want some instant data science that has more compute

Use jupyter-ai !



1. jupyter-lab

2. add an API key to OpenAI, Anthropic, DeepSeek or Qwen

3. use jupyter-ai

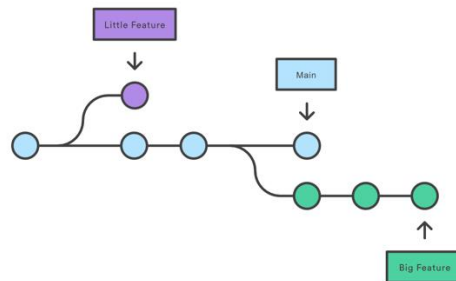4. edit your python notebook prolifically

jupyter-ai.readthedocs.io/en/latest/

# Ok, my code is getting complicated – multiple files, repository…

# Ok, my code is getting complicated – multiple files, repository…

Use aider.chat/ !

1. use git

2. new branch for every feature

3. use aider

4. (opt) in VSCode

# Ok, my code is getting complicated – multiple files, repository…

Use aider.chat/ !

1. use git

2. new branch for every feature

3. use aider

4. (opt) in VSCode
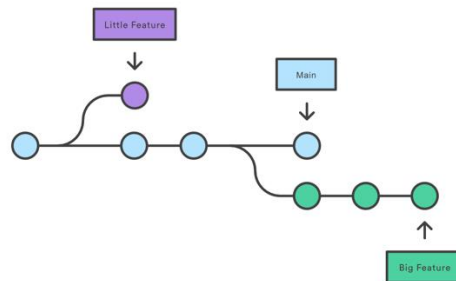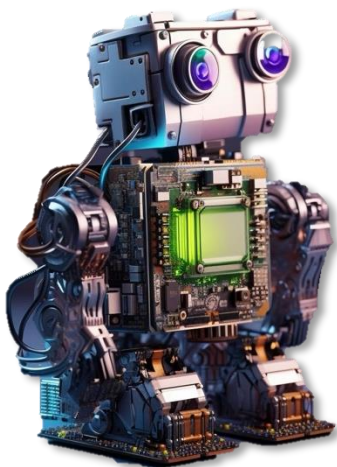
(it will commit every change, use `/undo` if it didn't work)

# You're trusting this guy

- He's doing his best

- But mistakes happen.

- Maybe you should pay him more

# Let's look in his brain



Percent completed correctly

claude-3-5-sonnet-20241022 · Qwen2.5-Coder-32B-Instruct (whole) · DeepSeek V2.5 · DeepSeek Chat V2 0628 (deprecated) · Dracarys2-72B-Instruct · gpt-4-0125-preview · ollama/Qwen2.5.1-Coder-7B-Instruct-GGUF:Q8_0-32k · o1-mini · Qwen2.5-Coder-7B-Instruct · gpt-3.5-turbo-1106 · Grok-2-mini · gemini-1.5-flash-exp-0827 · gpt-3.5-turbo-0613 · Codestral-22B-v0.1-Q4_K_M · yi-coder:9b-chat-q4_0 · Reflection-70B · Command R (08-24) · qwen1.5-110b-chat · ollama/codegeex4 · ollama/wojtek/opencodeinterpreter:6.7b

**aider code editing leaderboard**

# Let's look in his brain

```
app.py:

print("hello world")
```

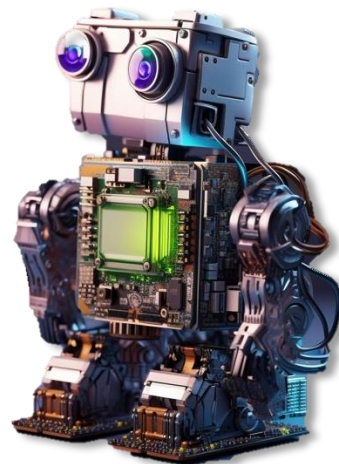aider.chat/docs/leaderboards/
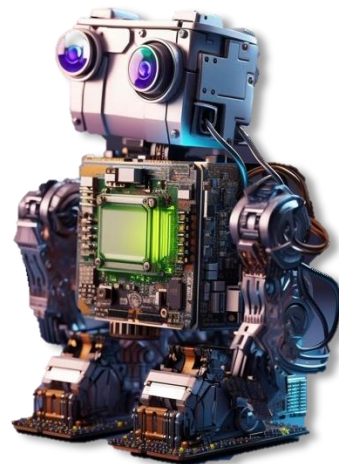**aider code editing leaderboard**

## Let's look in his brain

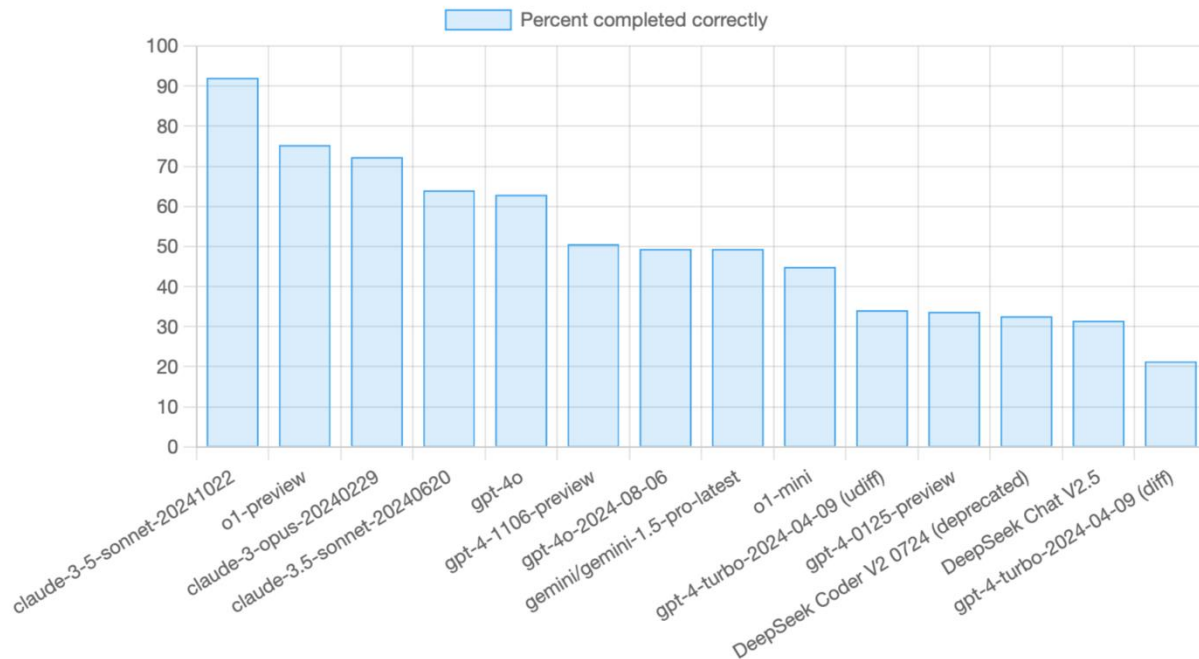**Refactor DatabaseSchemaEditor.quote_value**

Refactor the `quote_value` method in
the `DatabaseSchemaEditor` class to be a stand-alone,
top-level function. Name the new function `quote_value`,
exactly the same name as the existing method. Update any
existing `self.quote_value` calls to work with the
new `quote_value` function.

aider.chat/docs/leaderboards/
**aider code editing leaderboard**

# Let's look in his brain



Percent completed correctly chart with the following values:
- claude-3-5-sonnet-20241022: 92
- o1-preview: 75
- claude-3-opus-20240229: 72
- claude-3-5-sonnet-20240620: 64
- gpt-4o: 63
- gpt-4-1106-preview: 50
- gpt-4o-2024-08-06: 49
- gemini/gemini-1.5-pro-latest: 49
- o1-mini: 45
- gpt-4-turbo-2024-04-09 (udiff): 34
- gpt-4-0125-preview: 33
- DeepSeek Coder V2 0724 (deprecated): 32
- DeepSeek Chat V2.5: 31
- gpt-4-turbo-2024-04-09 (diff): 21

aider.chat/docs/leaderboards/
**aider refactoring leaderboard**

# Answer

Use
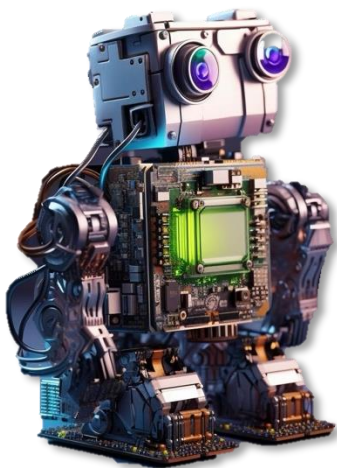`claude-3-5-sonnet-20241022`

Unless your stuff is private then use
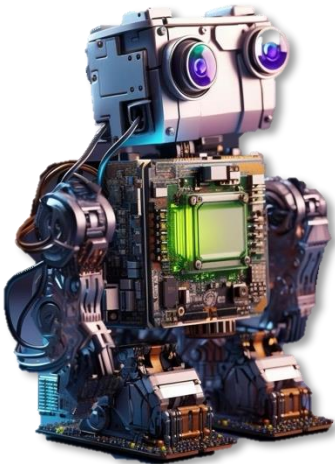`Qwen2.5-Coder-32B-Instruct`
or
`DeepSeek V2.5`

# Top tips

- Be as specific as you can about exactly what you want (or ask it to help you refine)

- make sure private stuff stays private!

- small incremental changes

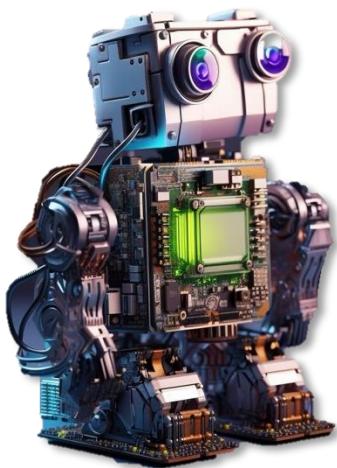- test as you go

# What's going to go wrong?



- Spaghetti code – not cohesive and difficult to maintain

- Functions with outputs they didn't need

- random imports that are different from what it said yesterday / in another file

- over complication

# What's not changing?



- We want to build our ideas, explore our data, understand the world

- We want tools that:
- helps us do this efficiently,
- minimize information overload
- don't miss critical parts,
- let us be the final judge
- help us make the important decisions, but not the irrelevant ones.

# Sydney Informatics Hub

**Sydney Informatics Hub** is a Core Research Facility within PVC-RI, enabling excellence in computational and data-driven research through advanced digital infrastructure, expert data consultancy and analytics training.

*Statistics*

*Data Science, AI & Software*

*Research Computing*

*Bioinformatics*

The University of Sydney     informatics.sydney.edu.au

# SIH Training

**64**

Different offerings, in person, online, and hybrid content in a variety of formats from webinars to interactive workshops.

Attendees at all career levels, from undergraduate students to senior professors, and representation from every Faculty and School.

Partnerships with national organisations like Australian BioCommons: biocommons.org.au/training-cooperative.

Find out more on our training calendar: sydney.edu.au/informatics-hub/training. Or stay up to date with our newsletter.

| Statistics | Data Science | Research Computing | Bioinformatics | Events |
|---|---|---|---|---|
| Fundamentals | Machine Learning | High Performance Computing | 'omics Techniques | Hacky Hour |
| Modelling | Visualisation | Cloud Computing | Reproducible Pipelines | Summer Schools |
| Specialist | Natural Language Processing | Containers | Data Analytics | Coding Challenges |
| | Geospatial Analysis | Workflows | National Compute Infrastructure | |

# Research examples with AI

# Open Australian Legal Corpus

Blog post: https://umarbutler.com/mapping-almost-every-law-regulation-and-case-in-australia/

# Extracting information from text with language models



NSW Supreme Court
Inheritance law judgments

| Name | James Nevill Dunn | ✅ |
|------|-------------------|----|
| Role in trial | Plaintiff | ✅ |
| Costs | $40,000 (p6) | ✅ |
| Natural person? | true | ✅ |
| Relationship to deceased | stepson (p1) | ✅ |
| Was dependent on deceased? | partly dependent (p32) | ✅ |
| Alleged misconduct? | false | ✅ |
| Estranged from deceased? | briefly estranged ~~due to a misunderstanding~~ (p43) | ✅ ❌ |
| Financial circumstances | receives a pension of $520 per fortnight, wife earns $800 per fortnight, owns a house with a $15,000 mortgage (p39-40) | ✅ |

OpenAI GPT-4o
(May '24)

Inheritance Case Law – Dr. Ben Chen
SIH: Mike Lynch, Xinwei Luo

https://www.caselaw.nsw.gov.au/decision/54a004453004262463c948bc

# Semantic search on private data

Digital Criminal Justice Project –
Dr. Carolyn McKay

SIH: Nathaniel Butterworth,
Gordon McDonald, Marius Mather
bit.ly/sih_search

Local language model
(not in the cloud)

Survey and Interview transcripts
⚠ Highly protected data

**Semantic Query**
food

⚙  Find

100%

before the sun had risen, the wife came and awakened the two children, saying, "Get up, you lazy bones; we are going into the forest to cut wood." Then she gave each of them a piece of bread, and said, "That is for dinner, and you must not eat it before then, for you will get no more." Grethel carried the bread under her apron, for Hansel had his pockets full of the flints. Then they set off all together on their way to the forest. When they had gone a little way Hansel stood still and looked back towards the house, and this he did again and again, till his father said to him, "Hansel, what are you looking at? take care not to forget your legs."

"O father," said Hansel, "Iam looking at my little white kitten, who is sitting up on the roof to bid me good-bye." - "You young fool," said the woman, "that is not your kitten, but the sunshine on the chimney-pot." Of course Hansel had not been looking at his kitten, but had been taking every now and then a flint from his pocket and dropping it on the road. When they reached the middle of the forest the father told the children to collect wood to make a fire to keep them, warm; and Hansel and Grethel gathered brushwood enough for a little mountain; and it was set on fire, and when the flame was burning quite high the wife said, "Now lie down by the fire and rest yourselves, you children, and we will go and cut wood; and when we are ready we will come and fetch you."

So Hansel and Grethel sat by the fire, and at noon they each ate their pieces of bread. They thought their father was in the wood all the time, as they seemed to hear the strokes of the axe: but
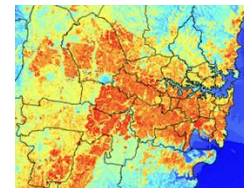
Survey   Workflows   Distributions   Data & Analysis

Edit question

Question type

Text entry

▼ Text type

Multiple lines

▼ Question behavior

Display logic

Skip logic

Default choices

JavaScript

The University of Sydney

test

XM   ≡   test

Tools   Saved at 10:57 AM   Draft

test

▼ Default Question Block

Q1

Q2
how long?

| | hours |
|---|---|
| thing1 | 3 |
| thing2 | 5 |
| thing3 | 17 |

# Understanding cities with computer vision

Prof. Ben Eggleton, NSW Smart Sensing, SIH: Henry Lydecker, Sahand Vahidnia, Thomas Mauch, Xinwei Luo
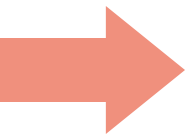
**Urban Heat ML Model**

**Geospatial Features**

Local Climate Zone = "2: Compact Midrise"

(49% class confidence)

**Compare neighborhoods**
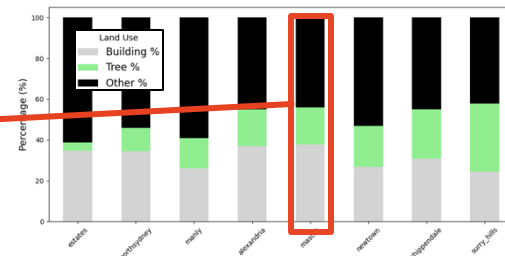
**Input Image**

bit.ly/aigis

**Computer Vision**

**Insights, ML, Analysis**

RGB Visual Light
Satellite Imagery
**Aerial Imagery** (Best)

Tree Patches

- 101 tree patches
- 15% coverage

Building Outlines

- 22 buildings
- 38% coverage
- building footprint area distribution

Building Size Distribution - mascot

# Speeding up microscopy data processing
## SIH: Sebastian Haan, Nathaniel Butterworth

bit.ly/sih-micro

## Key highlights

Sped up annotation process by up to 10x existing approaches

Full technology review of existing solutions with best practice solutions provided to clients

Strategic collaboration between SIH and SMM



Ground Truth

Model Prediction

Sydney Microscopy & Microanalysis collects immense amounts of data at the micro, nano, & atomic scales. We developed AI pipelines to speed up cell segmentation in electron microscopy images, automating 3D cell tomography.

The University of Sydney

SMM: Gerry Shami, Chad Moore, Andre Venne, Eleanor Kable, Prof. Filip Braet

# Australian Text Analytics Platform

*Text analysis tools for all researchers*



www.atap.edu.au

## Some tools on the platform:

**Document and corpus similarity tool**
- Compare differences between documents, e.g. to eliminate near-duplicates.

**Quotation tool and semantic tagging**
- Extract quotes from text e.g. news articles
- www.atap.edu.au/posts/quotation-tool/

**Discursis**
- An analysis and visualisation tool for conversational data
www.atap.edu.au/posts/discursis/

Linguistics: Prof. Monika Bendarek
SIH: Dr. Chao Sun, Hamish Croser, Jack Chan, Sony Jufri

# Generative X

**Drug discovery**
Design your own chemicals

**Biology**
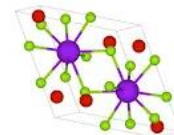Design your own protein
synthetic biology

**Manufacturing**
Design me a diffraction grating / micro
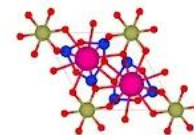resonator / widget with these properties…

**Material discovery**
AI mineral discovery and
autonomous material synthesis

Li$_4$MgGe$_2$S$_7$

KV$_3$SE$_3$

Rb$_2$HfSi$_3$O$_9$

Source: https://deepmind.google/discover/blog/millions-of-new-materials-discovered-with-deep-learning/

**Sydney Informatics Hub**
Dr Gordon McDonald, Informatics Team Lead
gordon.mcdonald@sydney.edu.au

**sydney.edu.au/informatics-hub**
sih.info@sydney.edu.au

THE UNIVERSITY OF
SYDNEY

Sydney
Informatics Hub

# Some Tools and Platforms

# Audio

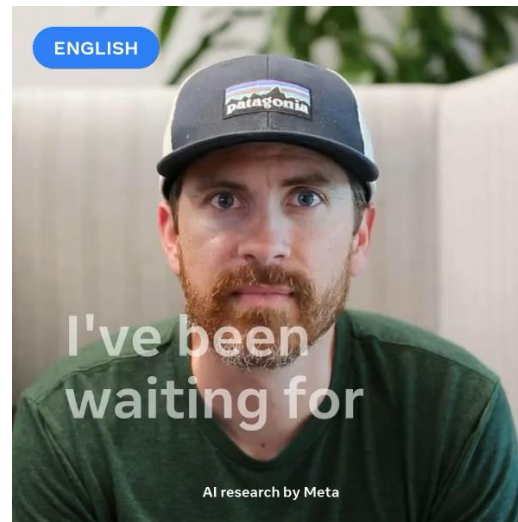**Translate audio** between different languages, keeping your tone of voice
https://seamless.metademolab.com/expressive

**Generate Speech**
https://elevenlabs.io

**Generate Music**
https://suno.com

# On-device video transcription

If you're on a mac
sindresorhus.com/aiko

Data and model runs in your browser – word level timestamps
huggingface.co/spaces/Xenova/whisper-word-level-timestamps

Data and model runs in your browser – speaker diarization
huggingface.co/spaces/Xenova/whisper-speaker-diarization

If you're ok with code and installs
github.com/ggerganov/whisper.cpp

(Not on device) Sharepoint automatic video transcription

# Image segmentation

Segment a human
huggingface.co/spaces/facebook/sapiens-seg

Segment images
sam2.metademolab.com

On-device - Data and model runs in your browser
huggingface.co/spaces/Xenova/segment-anything-web

Manual annotation
anylabeling.nrl.ai
roboflow.com

# Platforms & Tools

- [Huggingface](#) – Open & free tools for AI proof of concepts
- [Papers with code](#) – compare state of the art models' performance
- [Colab](#) – free Jupyter notebook interface from Google
- [Weights & Biases](#) – Model training and evaluation diagnostics
- [Roboflow](#) – Very easy to use dataset utilities
- [PaddlePaddle](#) - frameworks for AI use cases

# Literature Search

[Research Rabbit](#) – free for researchers, operates on donation model

[Elicit](#) free trial, then $10/month

[Rayyan](#) $8/month researcher, $4/month student

[Scite](#) Free trial then $13-$22/month

[SciSpace](#)

[Scopus-AI](#)

# Coding assistance

Free (ish)
Aider
Jupyter-AI
Open Interpreter
MetaGPT
Copilot

Paid
ChatGPT
Github Copilot

Best language models to use
with these tools now (10-09-24):

- claude-3.5-sonnet  (via API)

- gpt-4o (via API)

- DeepSeek 2.5 (via API)

- Local models? (none of them
  are quite good enough yet to be
  useful) ☹

# Agent building



**Agents**
to reduce
time to science

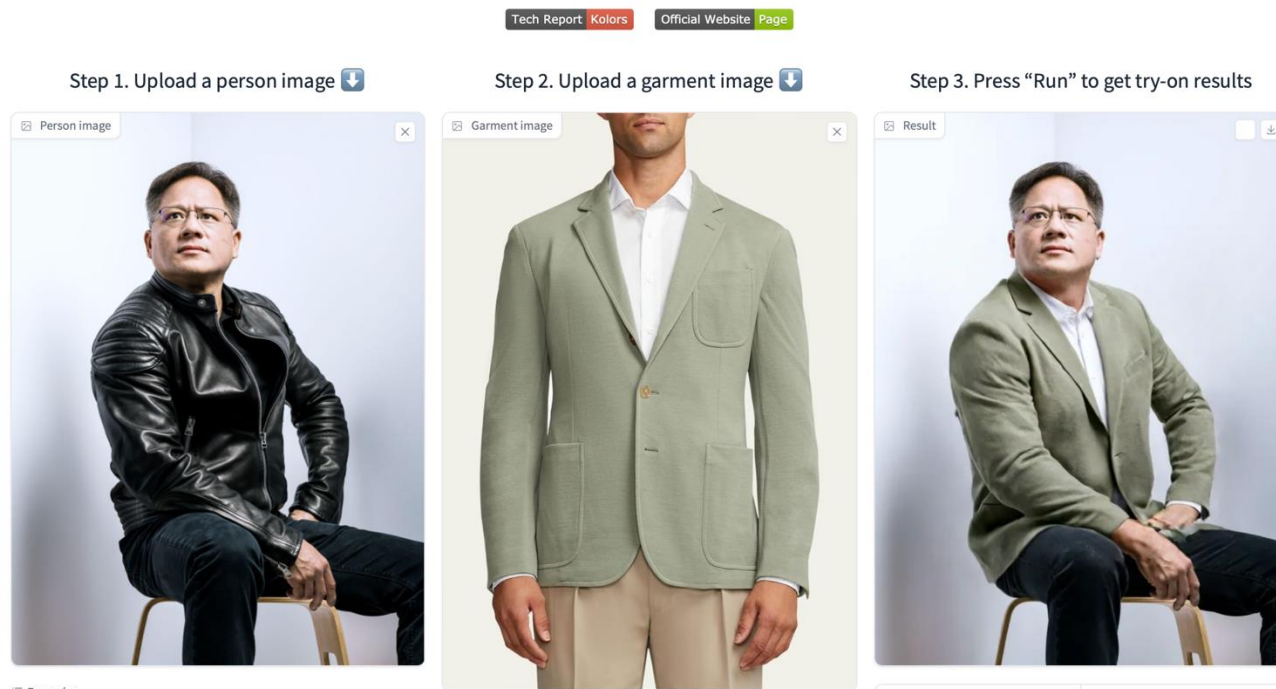Free (ish)
Cogniti
Flowise
MetaGPT

Paid
ChatGPT Custom GPTs
Azure Promptflow

# Try on virtual clothes!

To find more new models as they are released:
https://huggingface.co/spaces



**Kolors Virtual Try-On in the Wild**

Tech Report Kolors   Official Website Page

Step 1. Upload a person image ⬇    Step 2. Upload a garment image ⬇    Step 3. Press "Run" to get try-on results

https://huggingface.co/spaces/Kwai-Kolors/Kolors-Virtual-Try-On

**Sydney Informatics Hub**
Dr Gordon McDonald, Informatics Team Lead
gordon.mcdonald@sydney.edu.au

**sydney.edu.au/informatics-hub**
sih.info@sydney.edu.au

THE UNIVERSITY OF
SYDNEY
—
Sydney
Informatics Hub