

AI in Research

21st February 2024
Research Portfolio Researcher Spotlight

Dr Gordon McDonald
Dr Henry Lydecker

Sydney Informatics Hub

sydney.edu.au/informatics-hub



THE UNIVERSITY OF
SYDNEY

Sydney
Informatics Hub

CRICOS 00026A TEQSA PRV12057



We recognise and pay respect to the Elders and communities – past, present, and emerging – of the lands that the University of Sydney's campuses stand on. For thousands of years they have shared and exchanged knowledges across innumerable generations for the benefit of all.





Gordon McDonald

Six years building Quantum Sensors with atom interferometry at ANU
2016 – Now, University of Sydney, Sydney Informatics Hub.

In my time here we've had five thousand consults and projects in data science, AI, statistics, bioinformatics and high performance computing for academic researchers in all disciplines, and in collaboration with a myriad of national and international stakeholders. Together with my team, we've unleashed our tools and experience solving problems in conservation, medicine, law, geology, business, physics, education, genomics, archaeology and many more. Lately we're working out how to channel state-of-the-art AI models into research in a bunch of different ways which we'll tell you about in this talk.

Henry Lydecker

Redwood Forest Conservation Ecology
Coffee-ant Agroecology
Urban Ecology of Parasites
Zoonotic Disease Risk Modelling



Henry is a data scientist specialising in using artificial intelligence techniques including large language models and computer vision to solve problems in the real world. While his background is in ecology and public health research, he has worked on a range of projects from many different fields including agriculture, biology, ecology, geology, health, and law and is an expert in deep learning, machine learning, GIS, statistics, and software development. He leads a team of data scientists working to provide data science solutions to research problems, and to uplift the digital skills of the research community through training and outreach.



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



What is AI?

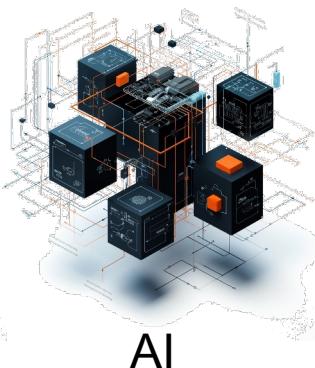
“A friendly koala cyborg”



what is the best way to peel a mandarin?



The University of Sydney



AI

has an internal representation
of meaning or content

[0, 1, .5, .2]



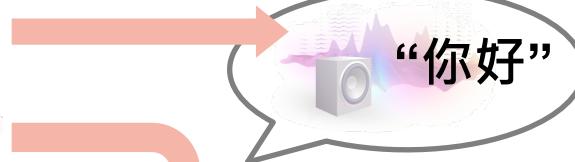
Generates...

Images, Video

Text



Peeling a mandarin, also known as a clementine or a tangerine, is generally a straightforward process. Here's a simple and efficient method to peel a mandarin:



Audio



3D structures
or molecules

Blessing

Amazing results in far less time/effort

Well-suited to poorly defined and ambiguous problems

Curse

Sometimes results will be unfeasible, wrong, ‘hallucinated’, biased etc.

You need a way to check the answer.

Some approaches are very data limited.

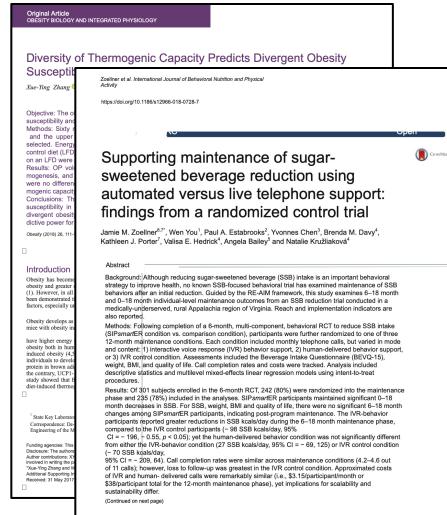
Modern Work

Information harvesting, integrating, and storytelling.

Large language models for literature review, synthesis, and text analysis.

Extracting information from text with language models

Academic papers on nutrition



“Who were the authors funded by?”

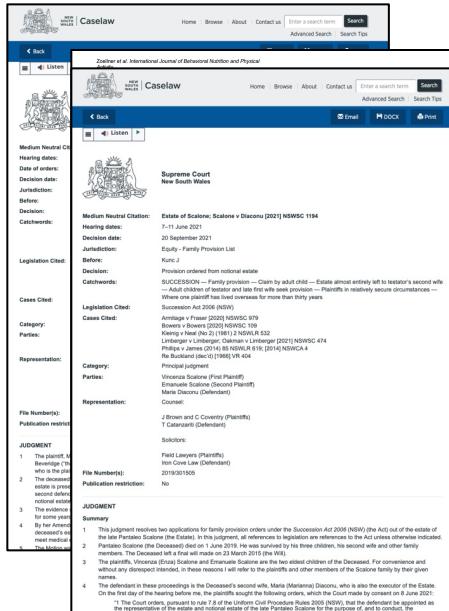
Extract

Funding body

Identifying biases in nutrition research - Dr. Fabian Held
SIH: Dr. Henry Lydecker, Dr. Joel Nothman, Dr. Chao Sun

Extracting information from text with language models

NSW Supreme Court Inheritance law judgments



The screenshot shows two pages of the NSW Caselaw website. The top page displays a search interface with fields for 'Hearing date', 'Date of orders', 'Decisions', 'Jurisdiction', 'Before', 'Decision', 'Catchwords', 'Legislation cited', 'Cases cited', 'Category', 'Parties', 'Representation', 'File Number(s)', and 'Publication restrict'. The bottom page shows a detailed case summary for 'Estate of Scalone; Scalone v Diaconi [2021] NSWSC 1154'. It includes sections for 'JUDGMENT', 'SUMMARY', and 'NOTABLE POINTS'. The judgment section lists several points, including the appointment of a personal representative and the distribution of assets.

Extracting Factors such as:

- Value of estate
- Relationship of parties to deceased
- Misconduct
- Drug or alcohol use

Extract

Structured data

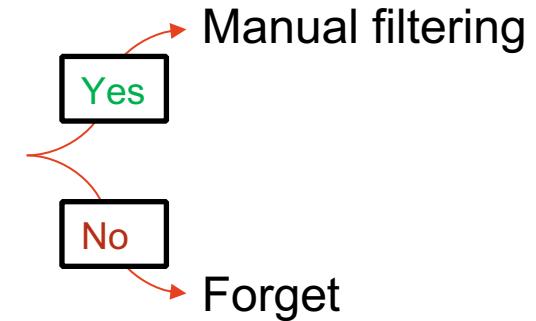
Inheritance Case Law – Dr. Ben Chen
SIH: Mike Lynch, Xinwei Luo

Extracting information from text with language models

PubMed articles

The screenshot shows the PubMed search results for a specific article. The search bar at the top contains the query "Improving de novo protein binder design with deep learning". Below the search bar, the article title is displayed: "Improving de novo protein binder design with deep learning". The authors listed are Nathaniel R Bennett, Brian Country, Jena Gorshenrik, Baekil Haeng, Azra Alen, Diana Velasquez, Ying Po Peng, Justus Daupens, Anilvijay Baek, Lance Stewart, Frank DiMato, Steven De Mureck, Savvas N Savvides, and David Baker. The journal is *Nat Commun.*, volume 6, issue 1471, published in May 2023. The DOI is 10.1038/s41467-023-38328-5. The article is freely available via PMC. The abstract section discusses the use of deep learning to design high affinity protein binding proteins. The full text is available through the nature portfolio link. The page includes standard navigation links like 'Save', 'Email', 'Send to', 'Display options', 'Cite', and 'Collections'.

Does this paper report
on an inherited trait?



OMIA Text Mining
Professor Frank Nicholas, Associate Professor Imke Tammem
SIH: Marius Mather, Sony Jufri, Joel Nothman, Di Lu

Searching for information

University of Sydney policies

The screenshot shows the 'LEAVE POLICY 2023' document. It includes the university logo, contact details for the Vice-Chancellor's delegate, and a section titled 'REASONABLE AND NON-ALLOWABLE EXPENSE PROCEDURES 2022'. This section contains details about the policy's issuance, amendment history, and signatory. A large red curved arrow points from this document towards the 'Policy Navigator' interface.

LEAVE POLICY 2023

The Vice-Chancellor's delegate
Dated:
Last Amended:
Signature:
Name:
Current Chancellor
CONT
1
2
3
4

REASONABLE AND NON-ALLOWABLE EXPENSE PROCEDURES 2022

Issued by: Chief Financial Officer
Dated: 27 October 2022
Last amended: 23 January 2024 (administrative amendments)
Signature:
Name: Wayne Andrews

1 Purpose and application

(1) These procedures are to give effect to the [Procurement Policy 2019](#) ("the policy"), and support compliance with probity, value for money and legislative requirements.

(2) These procedures apply to:

- (a) the University, students, staff and affiliates;
- (b) all foundations, centres, associations and institutes not separately incorporated;

Which policy document answers the question?

The screenshot shows the 'Policy Navigator' interface. A search query 'Can I take leave for jury service?' has been entered. The results page displays a card for 'JURY SERVICE' with a summary of the policy. Below the summary, there are two blue links: 'Policy Register Document URL: [Leave Policy 2021](#)' and 'Citations: [1. University of Sydney Enterprise Agreement 2023-2026.pdf](#), [2. Leave Policy 2021.pdf](#)'. The interface also includes tabs for 'Thought process', 'Supporting content', and 'Citation', and a sidebar with various policy sections and a 'Document Register' tab.

Draft version - use with healthy skepticism as AI-generated content may be incorrect. The authoritative source for University of Sydney policy is the official [policy register](#).

Can I take leave for jury service?

THE UNIVERSITY OF SYDNEY Policy Navigator

Thought process Supporting content Citation

Document Document Section Policy Register

337. Leave may be taken in one or more periods. Unused research or professional development leave does not accrue from year to year, and is not paid out on termination of employment.

JURY SERVICE

338. As a staff member who is summoned as a prospective juror must notify their Supervisor as soon as possible of the date when they are required to attend for jury service.

339. Upon providing proof of attendance, a staff member who is required to undertake jury duty will be granted:

- (a) paid leave for the period necessary for such service, in which case the staff member must travel to the University any fees (other than reimbursement of expenses) received for such jury service; or
- (b) leave without pay for the period necessary for such service, in which case the staff member will be entitled to receive any fees received for such jury service.

340. Where a staff member who takes paid leave fails to forward such fees to the University, the period of their absence for jury service will be treated as leave without pay, and any salary paid during that period of absence will be repaid to the University.

DEFENCE LEAVE

341. Staff are entitled to four weeks' paid leave per year to undertake Australian Defence Force (ADF) Reserve service training and operational duty.

342. An additional two weeks' paid leave may be taken in a staff member's first year of ADF Reserve service for attendance at recruitment and initial engagement training.

<https://policy-navigator.techlab.works>

Policy Navigator (alpha)

SIH: Henry Lydecker, Gordon McDonald

TechLab: Iqbal Chowdhury, Jim Cook

OGC: Lauren Myers, Deb Hook

Contextualizing information and storytelling

Papers, patents, PR, ...

Circulation Research
Volume 132, Issue 1, 6 January 2023, Pages 72-86
<https://doi.org/10.1161/CIRCRESAHA.122.321123>

ORIGINAL RESEARCH

Tropoelastin Improves Post-Infarct Cardiac Function

Meet the First Author, see p 5

Robert D. Hume, Shaan Kanagalingam, Tejas Deshmukh, Siqi Chen, Suzanne M. Mithieux, Fairoj N. Rasheed, Janee Roohani, Justyna Lu, Tapan Deo, Diana Graham, Zoe E. Clayton, E. Thomas, Anthony Weiss

Background: Myo... Following MI, necro... collagen, the extra... favorable properties introducing tropoelastin

Methods and Results: The method to administer to induced MI. Ex... PBS vehicle control assessments show (64.7±4.4% versus ms versus 31.1±5

Allergan to acquire University of Sydney spinoff Elastagen

9 February 2018

University spinoff Elastagen Pty Ltd has entered into a definitive agreement under which Allergan plc, a leading global biopharmaceutical company, has agreed to acquire the company.

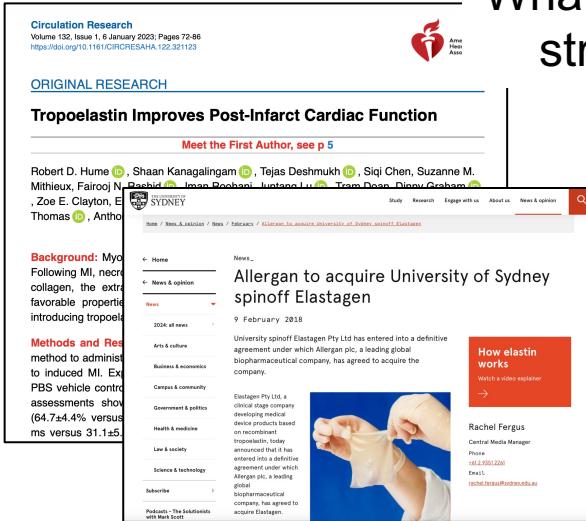
How elastin works

Rachel Fergus

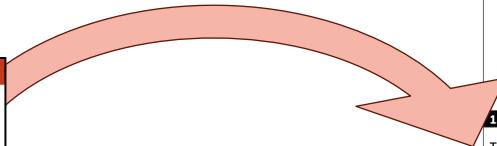
Central Media Manager

Phone: +61 2 9351 2220

Email: rachel.fergus@perkins.edu.au



What is the impact of this stream of research?



Research Impact Assessment Framework Case Study

Organisation Details		
Organisation name:	University of Sydney Charles Perkins Centre	
Title of case study:		Elastagen
Period when research was undertaken:		2013 – 2023
Period when the claimed impact occurred:		2013 – 2023
Names and roles of staff:	Staff name	Role
	Anthony Weiss	McCaughay Chair in Biochemistry, University of Sydney, NHMRC Leadership Fellow, Founder of Elastagen, and leader of Tissue Engineering & Regenerative Medicine in the Charles Perkins Centre, University of Sydney.
1. What is the problem your research seeks to address? Why is it significant? (250 word limit)		
<p>The research by Anthony Weiss and his team on Elastagen seeks to address the significant challenge of replicating the properties of human elastin, a key protein in the human body that provides elasticity and resilience to tissues such as skin, lungs, and arteries[1]. Human elastin's unique mechanical properties and biocompatibility make it difficult to mimic synthetically, which is a significant barrier in tissue engineering and regenerative medicine[2]. The significance of this research lies in its potential to revolutionize the field of biomedical engineering by providing a synthetic alternative to human elastin that can be used in a variety of medical applications, including wound healing, tissue repair, and the improvement of skin quality[1]. The development of Elastagen represents a breakthrough in the creation of biomaterials that can closely match the properties of natural human tissues, thereby enhancing the effectiveness of medical treatments and improving patient outcomes[2].</p>		

Research Impact Assessment Framework

SIH: Sebastian Haan, Nathaniel Butterworth, Gordon McDonald

FMH: Mona Shamshiri, Janine Richards, Prof. Robyn Ward

NSW Health

AI for Life Sciences

Genomics and Protein Design

AlphaFold

Dia-NN

RF Diffusion

AI laboratories

Enabling access to AlphaFold

SIH: Nate Butterworth, Chris Le, Nandan Deshpande

FMH: Jake Chen

Key Highlights



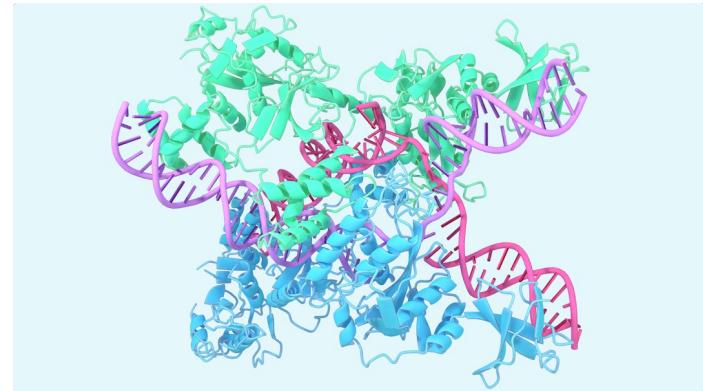
Deployed on Artemis
HPC, NCI Gadi, Galaxy
Australia, Ronin, &
Google Colab.



Training materials publicly
available on GitHub &
YouTube.



USYD researchers now
able to easily apply
AlphaFold in their work.



Google Deepmind's AlphaFold2 has revolutionized the way that researchers study protein structures.

We've provided USYD researchers many ways to access and use AlphaFold, and developed training courses & materials to make it easier for a diverse range of researchers to use this cutting-edge tool in their work.

<https://sydney-informatics-hub.github.io/training.alphafold/>

Scaling Proteomics with Dia-NN

SIH: Cali Willet, Tracy Chew, Nathaniel Butterworth, Georgie Samaha

Sydney Mass Spec: Ben Crossett, A/Prof Carsten Schmitz-Peiffer (SOLES)

Key Highlights



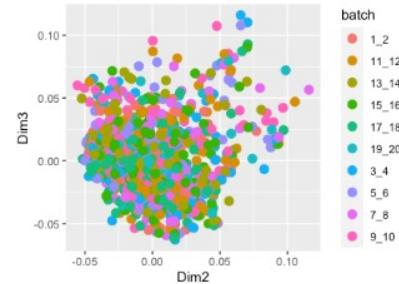
Deployed state-of-the-art tool for identifying proteins in mass spec data on NCI GADI.



Eliminated data processing bottleneck, speeding up by orders of magnitude (46 days down to 18 hours).



Strategic collaboration between SIH, Sydney Mass Spec, Australian Biocommons, and the CPC.



Dia-NN is the leading tool to identify peptides and proteins in mass spectrometry data. We built a scalable Dia-NN workflow, dramatically speeding up processing and solving a major bottleneck in large scale proteomics studies.

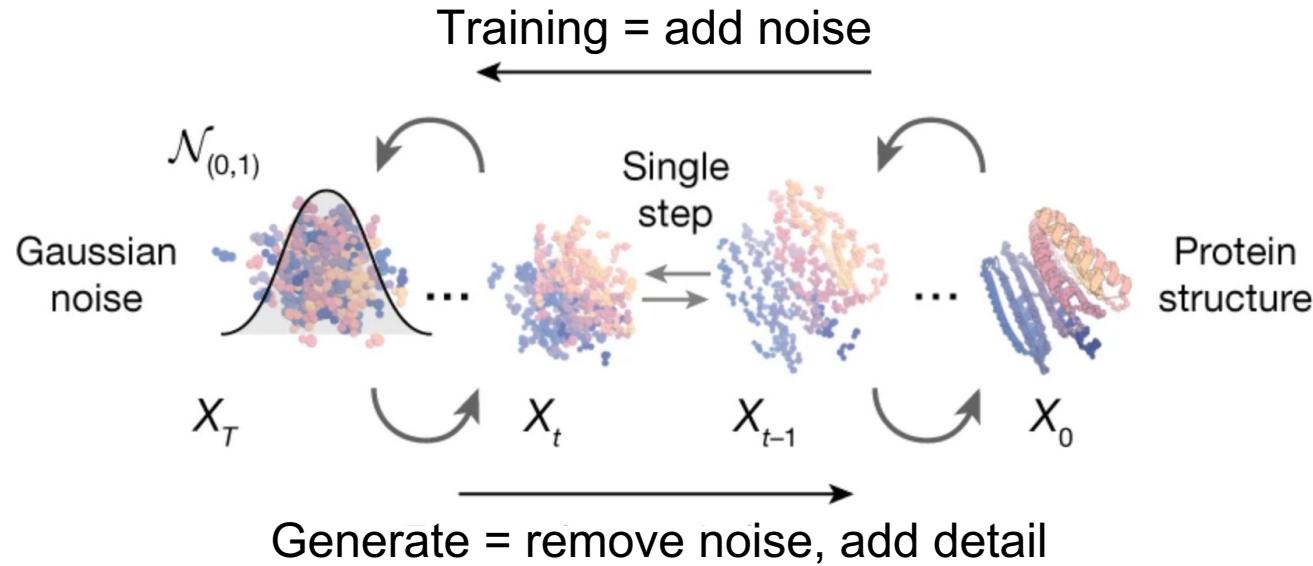
Diffusion models

Training = add noise

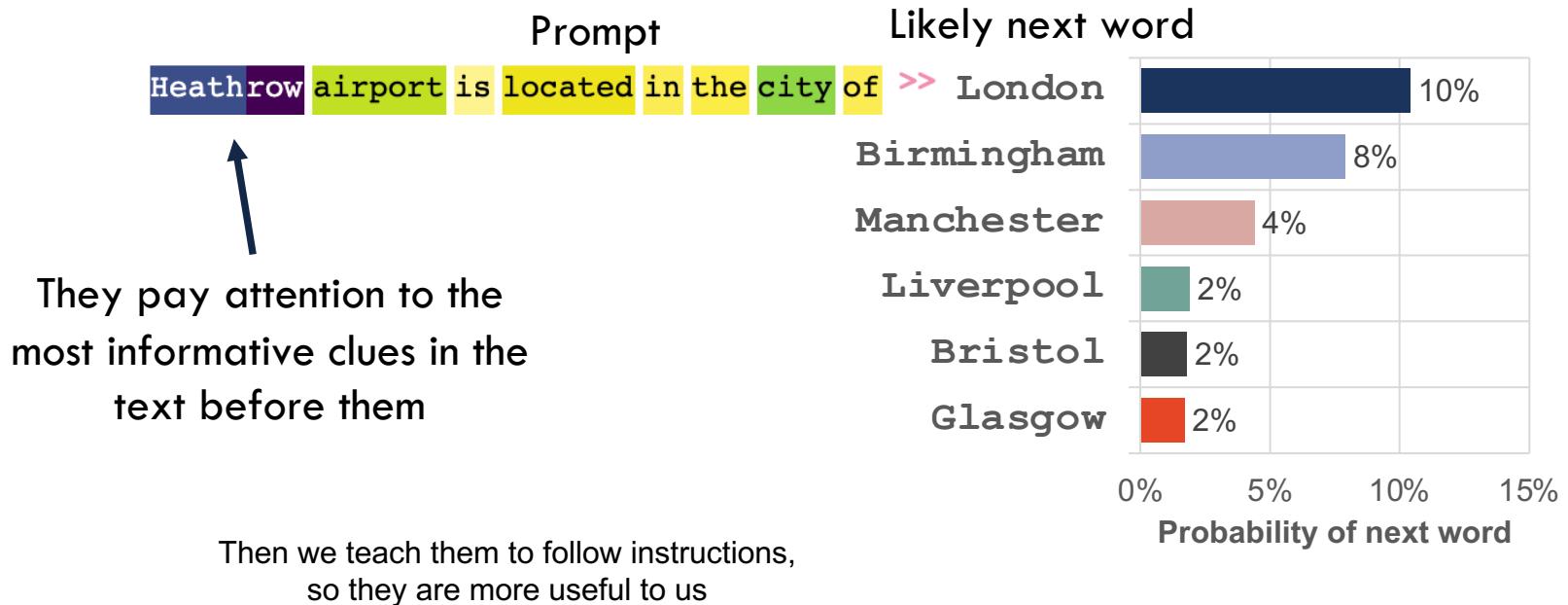


Generate = remove noise, add detail

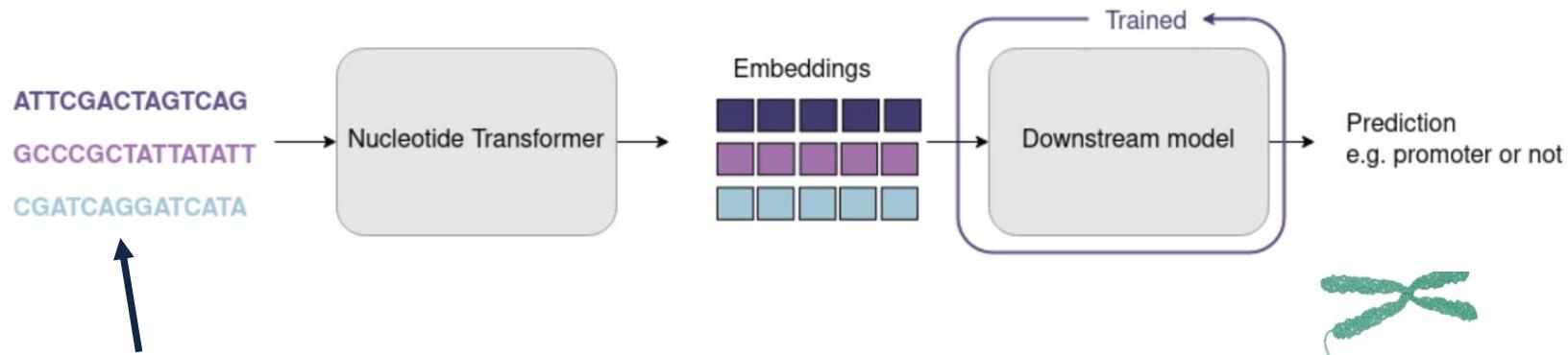
Diffusion models



Language models (like ChatGPT)

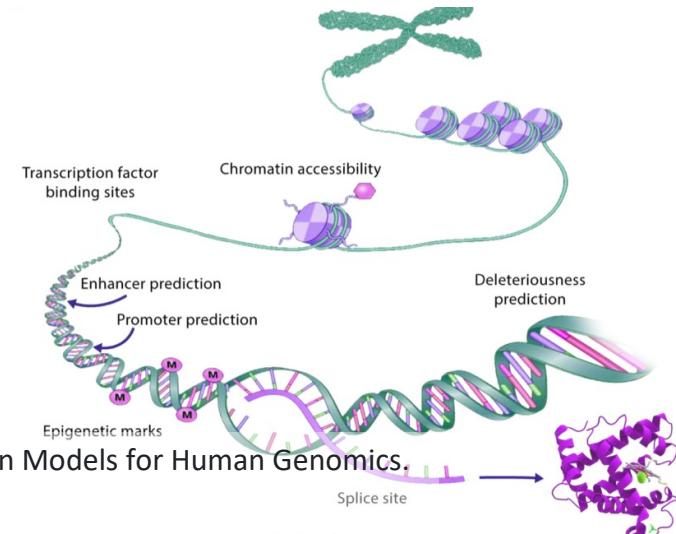


Genomics models



They pay attention to the most informative clues in the text before them

Then we fine tune them, so they are more useful to us for particular tasks



Agents



LLMs can simulate individual roles in a project

- Project manager
- Research assistant
- Professor
- Software engineer
- Statistician
- Subject matter expert.



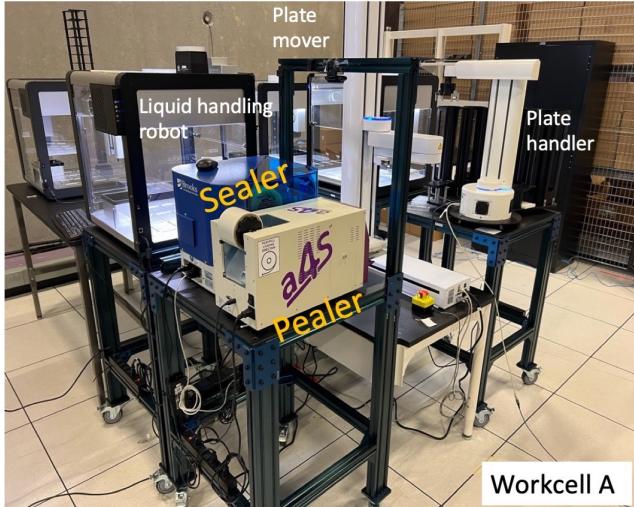
Make them talk to each other, and you can
create a team of specialist AI agents
working together towards your goal.



Get results in minutes instead of days!



Agent-based lab automation



Self Driving Laboratories @ Argonne
<https://github.com/AD-SDL>

The University of Sydney



A-Lab, a facility at Berkeley Lab where artificial intelligence guides robots in making new materials created by GenAI

Szymanski, et al (2023). *An autonomous laboratory for the accelerated synthesis of novel materials*, Nature
Merchant, et al (2023). *Scaling deep learning for materials discovery*, Nature

Picture This

AI in Geospatial, Microscopy, Medical, Agriculture, & Conservation imagery

Image Classification
Object Detection
Instance Segmentation
Multi-Modal Models

“We have a lot of images, and we want a way to automate analyzing them”

Problem: Ecologist has 100,000 images and wants to know what is in them

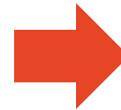
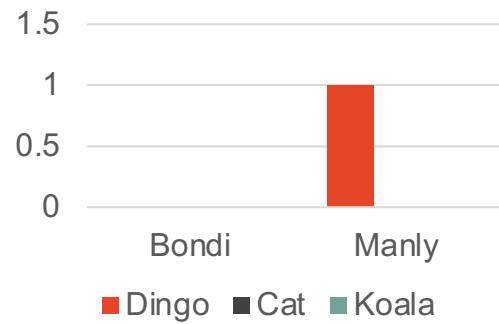


Photo	Species	Count	Site
1	NA	NA	Manly
2	Dingo	1	Manly
3	NA	NA	Manly
4	NA	NA	Manly
5	NA	NA	Manly
6	NA	NA	Manly
7	NA	NA	Manly

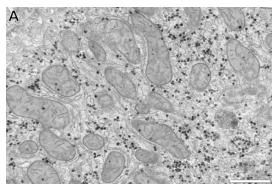


Dr Aaron Greenville, SOLES

The University of Sydney



Computer Vision



Model

What is this image of?

Is there anything in it?

Where are the things in it?

Detecting wildlife with NSW National Parks

SIH: Henry Lydecker, Nathaniel Butterworth, Gordon McDonald



Key Highlights



Fine-tuned family of YOLOv5 models on **2.16 million** camera trap images.



Detect and identify up to **72 different species**, with **94% precision** and **95% recall** on the top 33 most common.



GPU-accelerated prediction pipeline **1371x faster** than manual labeling.

Ecological researchers like Dr. Aaron Greenville, conservation agencies, and governments make extensive use of camera traps to monitor wildlife. Processing this data requires an immense amount of time.

We implemented an AI pipeline for automatically detecting & identifying animals in camera trap images, dramatically speeding up speed of data processing.

Speeding up microscopy data processing

SIH: Sebastian Haan, Nathaniel Butterworth

Key highlights



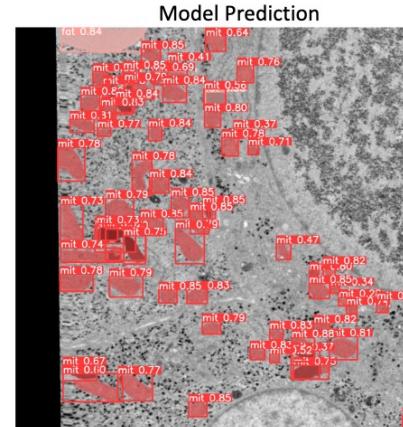
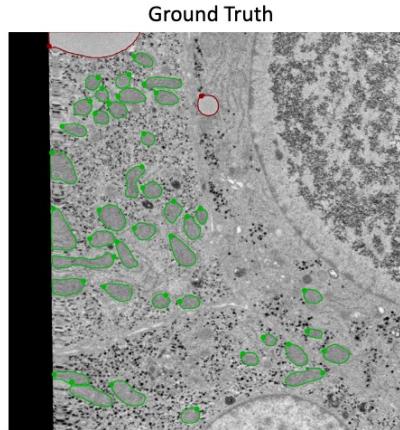
Sped up annotation process by up to 10x existing approaches



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



Strategic collaboration between SIH and SMM



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

Using urban morphology to model thermal comfort

SIH: Henry Lydecker, Sahand Vahidnia, Thomas Mauch, Xinwei Luo

Key Highlights



Open-source Python tools for spatial imagery processing, annotation, and analysis using AI.



Model training dataset:
55k tree patches, 75k buildings, 3,900 aerial images.



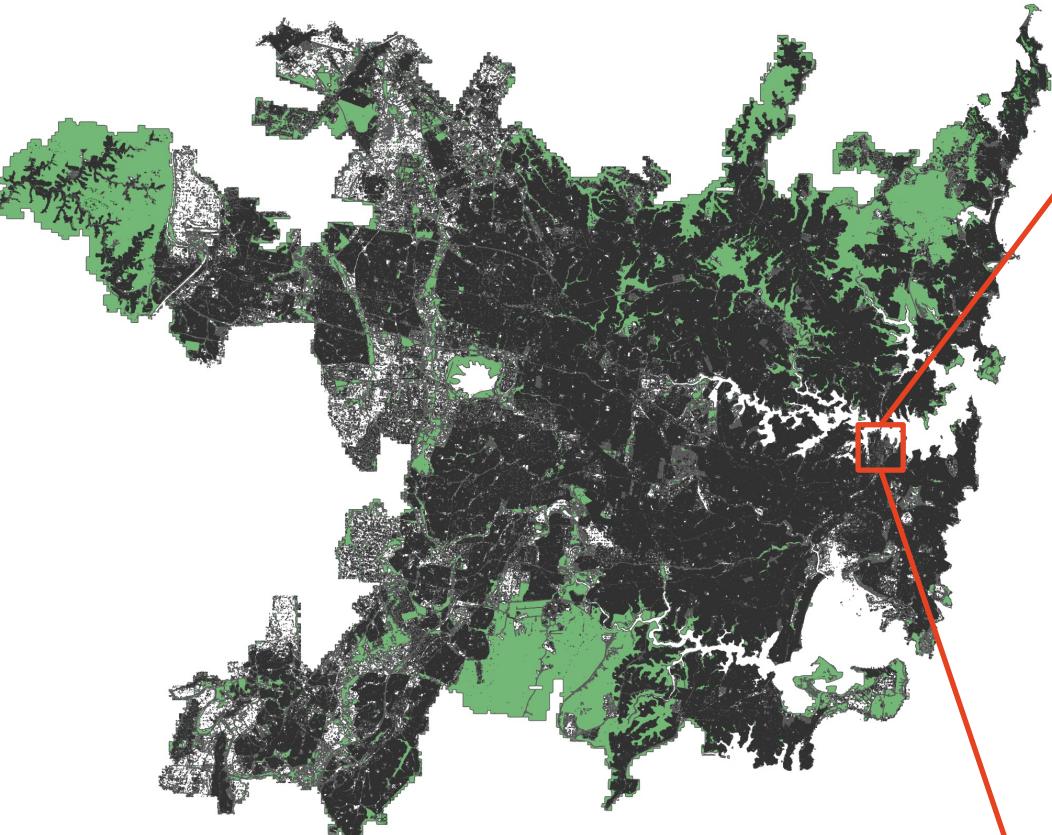
Predicted Dataset: 2 million tree patches, 1 million buildings, 37k aerial images over 12k km². Freely available in GeoParquet format.



Professor Ben Eggleton's group are studying how heat waves are harming cities, increasing water demand and health risks.

We developed AI models for extracting key variables for understanding human thermal comfort from spatial imagery, paving the way for future solutions to support urban cooling and water demand management.

Greater Sydney Dataset



>2 million tree patches
1 million buildings



THE UNIVERSITY OF
SYDNEY
Sydney
Informatics Hub

Understanding cities with computer vision

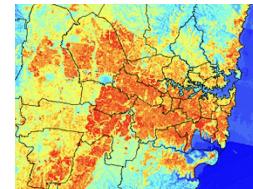
Prof. Ben Eggleton, NSW Smart Sensing, SIH: Henry Lydecker, Sahand Vahidnia, Thomas Mauch, Xinwei Luo



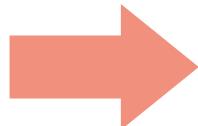
Geospatial Features

Local Climate Zone = “2: Compact Midrise”

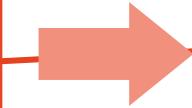
(49% class confidence)



Input Image



Computer Vision



Insights, ML, Analysis

RGB Visual Light
Satellite Imagery
Aerial Imagery (Best)

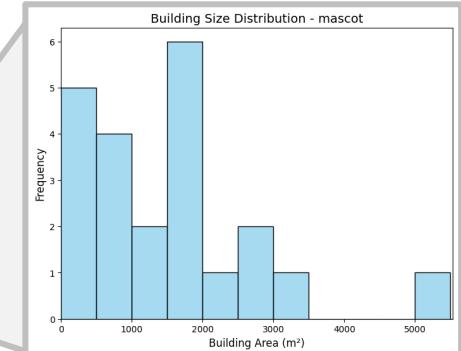
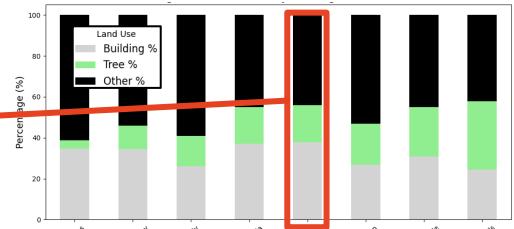
Tree Patches

- 101 tree patches
- 15% coverage

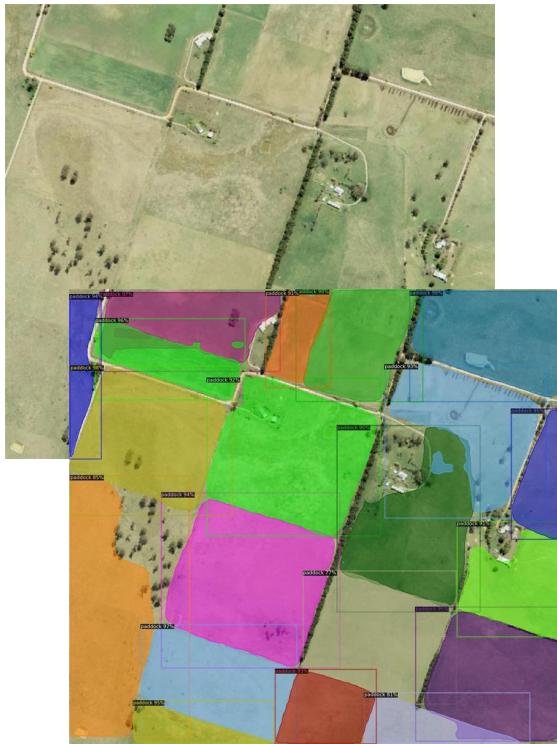
Building Outlines

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

Compare neighborhoods



Agricultural land use models, urban ecology



This tool will use the produced data to compare the soil moisture of a farm against historic data. Please select the desired comparison method and dates to make the comparison as in section A. Then choose the visualisation in section B to see the results.

A

Rmaapl

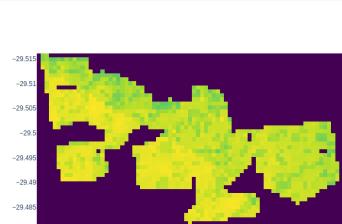
Select soil layer: SMI Select the historic years to compare against: 2023-02-01 → 2023-02-28

Select window aggregation method: std Select historic aggregation method: std

Generate Images

B

Select visualisation name: Historic_2023_02_01_2023_02_28_sm1_w_std_h_std_y_5 Select your palette: viridis



Recommended Tools and Platforms

Computer Vision Models

- [YOLOv5](#) – models for classification, object detection, and segmentation
- [Detectron2](#) – powerful tools, using vision transformers
- [Segment Anything Model](#) – general purpose zero-shot panoptic segmentation
- Grounding DINO – make your models multi-modal using text

Platforms & Tools

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

Literature Search

Research Rabbit – free for researchers, operates on donation model

Consensus Free basic usage, \$7/month for pro features

Paper Digest Literature Review Free to try but \$7-\$11/month for pro features

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)

Jupyter-AI

Open Interpreter

Tabby

MetaGPT

Paid

ChatGPT code interpreter

Github Copilot

Agent building

Free (ish)

Cogniti

Flowise

MetaGPT

Paid

ChatGPT Custom GPTs

Azure Promptflow



Sydney Informatics Hub

Dr Gordon McDonald, Informatics Team Lead
Dr Henry Lydecker, Data Scientist
sih.info@sydney.edu.au
first.last@sydney.edu.au

sydney.edu.au/informatics-hub



THE UNIVERSITY OF
SYDNEY

—
Sydney
Informatics Hub

