# White Paper - Decadal Plan Working Group 2.3 - Data and Computing

Chairs: Minh Huynh (CSIRO) and Chris Power (ICRAR/UWA)

Sub-Working Group Leads: JC Guzman (SKAO), Minh Huynh (CSIRO), Mark Krumholz (ANU), Simon O'Toole (AAO/Macquarie), Chris Power (ICRAR/UWA)

## Executive Summary

The Decadal Plan community survey results reveal that more than 50% of the Australian astronomical community make use of data and computing infrastructure and services as part of their research. Engagement with the community revealed some recurring themes.

- Data platforms are now established as a critical component of how researchers engage with their data, and are regarded as essential by observational astronomers interacting with modern observatories. The All-Sky Virtual Observatory (ASVO), including the likes of the AAO Data Central science platform, the CSIRO ASKAP Science Data Archive (CASDA), and the MWA ASVO, is regarded as an outstanding investment. There is strong support for expansion of funding to support increased capacity and coverage of data platforms and archives, further developments of their capabilities, and their integration into a common astronomy science platform.
- The mid-scale HPC facilities OzStar and NT at Swinburne University of Technology, open access to all astronomers in Australia upon request, and provide critical computing infrastructure for the majority of the community; continued support and investment in them should be prioritised. The peak HPC facilities at National Computational Infrastructure and Pawsey Supercomputing Centre serve the needs of the largest users, primarily in theoretical astrophysical simulations and radio astronomy. However, the amount of compute time currently available that is open access is insufficient and is unable to meet the needs of the theoretical astrophysics community, while there have been challenges in meeting the operational requirements of the Square Kilometre Array (SKA) precursors, ASKAP and MWA.
- The importance of a professional software base developed according to best practice is widely recognised. The Astronomy Data and Computing Services (ADACS) has emerged as one of the key successes during the previous Decadal Plan period. Ongoing support for an expanded and sustainable ADACS is regarded as a priority, as well as more general investment in professional data and software engineering in the development of software and pipelines built into the costing of projects.
- Community-wide training to increase data and computing competency is important; this should be via a nationally coordinated training and accreditation scheme, incorporated into PhDs but available community wide, but can also happen via communities of practice to bring together researchers within disciplines.
- Clear pathways for career development - and long-term security - for the data, software, and platform specialists that underpin the astronomical community's data and computing needs is critical.

## Recommendations

- **Software**
  - Invest in Australia's data reduction pipeline capabilities - including scalability and automation - and improve funding for development and long-term support of pipelines.
  - Invest in widely-used mission critical software with explicit funding to ensure software security; this should be a condition for funding of new telescopes and instruments.
  - Establish a mechanism for long-term astronomical software sustainability to ensure the impact and legacy of these codebases.
  - Ensure the continued funding of ADACS - preferably at an expanded level.
  - Support software migration from CPU- to GPU-based architectures, which increasingly dominate the top end of HPC hardware, and to maximise software and workflow performance on heterogeneous architectures in current and future HPC facilities.

- ○ Support research software engineering in general (via professional bodies such as [Research Software Engineers Society of Australia and New Zealand](#))
- **Standards**
  - ○ Ensure that all Australian-hosted data are [IVOA standards compliant](#), adheres to [FAIR Principles](#), and is accessible via data portals that adhere to international standards, such as [Web Content Accessibility Guidelines 2.1.](#)
  - ○ Ensure that all Australian-developed software is [IVOA](#) standards compliant; adheres to the [software FAIR Principles](#); and is [Open Source](#).
- **Platforms**
  - ○ Support key national research data platforms to provide long-term availability of data, via ongoing funding and regular capital refresh.
  - ○ Invest in an Australian astronomy-focussed science platform, allowing interactive analysis of remote datasets that fuses together multi-wavelength, multi-messenger, and theory/simulations data.
- **Open-access HPC**
  - ○ The astronomical community should continue to support [OzSTAR](#) and [NT](#) and the programs initiated by [AAL](#) to provide computing time on [NCI](#). [AAL](#) should consider extending this access program to the [Pawsey Supercomputing Centre](#), and should investigate other options including commercial providers and overseas facilities.
  - ○ Australia should move more of its computing resources into an open, merit-based allocation scheme, with the goal of reaching approximately 30-40% open access – in line with peer countries – over the decade.
  - ○ There should be enough HPC time available for it to be possible to carry out large programs comparable to those that are possible in peer nations - a target is 30% of a top-100 machine is recommended.
- **Supporting the SKA and its precursors**
  - ○ The continued support of [ASKAP](#) and [MWA](#) by stable and reliable HPC systems is critical for the real-time, high throughput compute required by the processing and analysis workflows to support surveys in the coming decade.
  - ○ Opportunities to refresh operations-critical compute platforms for high-data-rate telescopes should consider the unique requirements around sustained data throughput and the emphasis this places on storage, filesystems and isolation from general users to ensure reliability.
  - ○ Support [AusSRC](#) to provide resources to support radio astronomy research, in addition to resources available at [Pawsey](#) and [NCI](#), ensuring that the resources keep pace with the growing size of the SKA arrays and are in place in a timely fashion to avoid lost opportunities.
- **Training**
  - ○ Establish a nationally coordinated training and accreditation scheme for data & computing competency, incorporated into PhDs but accessible to early career researchers and the community as a whole.
- **Working Practices and Careers**
  - ○ Stable and predictable pathways for career development and long-term security are essential for the data, software, and platform specialists that underpin the astronomical community's data and computing needs.
  - ○ Complex data and software projects should be explicit collaborations between domain specialist researchers and specialist data, software, and platform engineers who can build the codebases and workflows adhering to best practice.
  - ○ Establish communities of practice - groups of researchers within specific sub-disciplines (e.g. AI/ML) that collectively manage their shared software and data needs, working with teams of dedicated infrastructure specialists.
  - ○ Effective adoption of advancing novel technologies such as XR/VR/AR to enable new and innovative pathways to scientific outcomes alongside efficient distributed collaboration methods

# Detailed Commentary

In the following sections, we summarise the main findings of the sub-working group reports, providing detailed commentary, highlighting specific topics of importance, and summarising the key recommendations.

## 1. Community Perspectives on Data

**General Comments:** Data is one of the fundamental and necessary prerequisites for Australian astronomy's research programme. It encompasses raw or reduced/processed data from telescopes or equivalent facilities; simulated data of physical systems, either astrophysical or instrumental; the results of numerical computation or theoretical modelling; and derived data products from analyses of any of these. It is generated, handled, and stored; it may be reduced and processed; it requires analysis and interpretation; and it needs to be findable and accessible by the community. As a result, our community's approach to data requires the expertise of data engineers and data scientists, and if we are to extract its maximum value, it requires data coordination between data centres, data and software engineering teams, funding agencies, and industry.

The previous Decadal Plan discussed data in the context of data-intensive research. It mentioned the "missing layer in the era of data intensive research" (p.54) with particular reference to the All-Sky Virtual Observatory (ASVO), which has played a key role in linking Australian data engineering and science with international counterparts. However, there were no data-specific infrastructure recommendations or goals; there were data-related recommendations for better support for HPC and related software, and these areas have seen increased funding and support over the last decade. While the Mid-term Review briefly discussed the International Virtual Observatory Alliance (IVOA), which has played a key role in data intensive projects such as Gaia and in planning for the Vera Rubin Observatory (VRO) and SKA, it has been noted that securing funding and support for ASVO has been challenging over the last decade. However, the Mid-Term Review made the key recommendation,

> *Astronomy should build an astronomical data fabric that links high-performance resources through appropriate data middleware and networks to create new opportunities for discovery by Australian researchers based on data flowing from telescopes like SkyMapper, ASKAP and the MWA.*

This vision is being realised as we approach the end of the previous Decadal Plan's term, through the Data Central Science Platform, the CSIRO ASKAP Science Data Archive (CASDA), the MWA ASVO, the Gravitational Wave Data Centre (GWDC), and the Australian SKA Regional Centre. While a facility such as the MWA ASVO is rightly regarded as a spectacular success, it has faced periods of uncertainty around its funding and operated on a relatively constrained budget. If the vision set out in the Mid-Term Review is to be fully realised, sustained and stable funding and fully committed support will be required, which can only enhance the impact of the various ASVO facilities and AusSRC. Moreover, building stronger links between multi-wavelength and multi-messenger data platforms will be critical to open up new discovery spaces and maintain high research impact into the future.

**Research Data Engineering:** The rise of research software engineering over the last decade is mirrored in the growing importance of research data engineering. Research data engineering encompasses data handling, storage, processing and reduction, as well as processes such as cleaning and quality assurance - all the stages in the data life cycle before it is ready to be analysed by a researcher. Data should adhere to the FAIR principles, which were formulated in a 2016 Nature article - data should be Findable, Accessible, Interoperable and Reusable - although the basic concepts are older.

Access to specialised data engineering support often depends on the size of the institution or research team, and the relative maturity of the field. For example, large observatories and related facilities will provide many of these services; large numerical simulations teams carefully curate their data products before making them available for more general use; and more established fields, such as optical/IR and radio astronomy, have had access to greater data engineering resources available in recent times. In contrast, smaller institutions and research teams, and younger fields, such as gravitational wave astronomy, rely almost entirely on groups of

researchers to carry out their data engineering. Providing more uniform access to research data engineering support via investment and training is becoming more pressing, as the diversity and scale of datasets increase.

**FAIR Hosting of Astronomical Data:** Data handling and storage is the process of acquiring data in a (sometimes) remote location and moving it to where it can be stored and accessed for processing, reduction, and subsequent analysis. Australian astronomy has led the way in making data available to the community. Indeed, Australia was a founding member of the IVOA, which has the goal of developing standards and protocols for data interoperability and accessibility.

There has been significant investment by the National Collaborative Research Infrastructure Strategy (NCRIS) into astronomical data infrastructure, including the Data Central Science Platform, the Murchison Widefield Array archive, the SkyMapper archive, CASDA, the Theoretical Astrophysical Observatory (TAO) and GWDC. Each of these facilities hosts hundreds of terabytes to tens of petabytes worth of data. There are four centres that support them: the National Computational Infrastructure (NCI – SkyMapper and a small part of Data Central), the Pawsey Supercomputing Centre (MWA and CASDA), OzStar (TAO and GWDC) and Data Central's facility at AAO (Macquarie University). A collaboration between OzStar and Data Central will host catalogue data from the Legacy Survey of Space and Time (LSST) when that survey data becomes available.

Accessing these data is facilitated by implementation of IVOA standards, which are now in place at CASDA, Data Central, the MWA ASVO archive and SkyMapper; all Australian-hosted data should be made IVOA standards compliant. This will ensure that these datasets are (for the most part) adherent to FAIR principles, and improve data accessibility.

Services such as the Table Access Protocol (TAP), Simple Image Access (SIA), and Simple Spectral Access (SSA) allow astronomers to access a huge amount of data in a standard way, through tools such as TOPCAT and Aladin. Each of the systems listed also allow querying via a web portal. There should be widespread adoption of Web Content Accessibility Guidelines 2.1 and consideration of e.g. colour vision deficiency when designing and developing user interfaces, to enhance accessibility and to improve user experience.

There are challenges ahead. The SKA; optical/infrared facilities like VRO and Extremely Large Telescopes; space telescopes like Nancy Grace Roman; and next generation numerical simulations, will generate enormous data volumes, and so data storage will need to be distributed. This increases the value of the IVOA that provides standardised ways to query and process data in a distributed way. For example, Data Central's Data Aggregation Service provides an interface that allows users to bring together data from across the ASVO in a single interface; the ASVO has Single Sign On capabilities across each of its five nodes. Such data storage will need to be long-term and (relatively) readily accessible, which is especially important as data volumes grow but the number of researchers remains roughly constant.

**Processing Astronomical Data:** There are two aspects to data reduction and processing: the development of pipelines and the running of those pipelines. This is one of the most important aspects of research data engineering; without robust, managed, and well-maintained data pipelines that have been developed according to best practice, our community's ability to extract value from its data is limited, especially in the coming era of enormous datasets, and Australian astronomy's international standing is at stake.

How development of data reduction and processing pipelines are supported and managed has changed over the last decade. At the time of writing of the previous decadal plan, data reduction pipelines for Australian facilities were largely developed and maintained by nationally supported observatories, such as the AAO and CSIRO ATNF. Since then, there has been significant investment in radio astronomy in the development of YandaSoft, the real-time calibration and imaging pipeline for the ASKAP radio telescope. This contrasts with optical astronomy; the 2dFdr pipeline for the 2dF and AAOmega spectrographs on the AAT is maintained at a minimal level, while other instruments on the AAT currently do not have stable and easily accessible pipelines, and so difficult to maintain code that can only be run by a handful of researchers on specific computers. Future instrument build projects should adhere to what might be termed the ESO model: invest in the

development of properly engineered data reduction pipelines that are robust, tested and maintainable. This will mean faster and greater impact for any and all instruments that Australia builds in future. Legacy data reduction pipelines for existing instruments should be supported as part of an Australian-based software development and maintenance fund.

The growth in data volumes and the increasing complexity of reduction and processing requirements over the last decade means that the computing requirements of pipelines in most areas of astronomy now exceed high-end laptops or desktops, and must be carried on computing clusters or high performance computing facilities. This has been the case in radio astronomy and astrophysical simulations for some time, but is increasingly the case in optical and infrared astronomy following recent technological advances in instrumentation (e.g. advent of large scale integral field spectroscopy; short cadence, large field imaging). As a result, the computational and technical requirements of modern data reduction and processing pipelines require expert level users to run them, trained in high performance computing and data but with the domain expertise to understand the data; this places constraints on the size of science team that could carry out this kind of data processing, unless there is additional support from organisations such as Data Central or AusSRC. It's also worth noting that mid-scale computing facilities, such as university maintained computing clusters or Swinburne University of Technology's OzStar/Ngarrgu Tindebeek supercomputer, are well matched to the requirements of pipelines in optical/infrared, and gravitational wave astronomy, while NCI and Pawsey continue to play a vital role for radio astronomy and large-scale astrophysical simulations.

Australia has an excellent opportunity to build on its existing capabilities for pipeline automation at scale for users of existing and future large and complex instruments. There is a gap between laptop/desktop computation and HPC that urgently needs to be filled, to support the next generation of large and complex instruments. Groups at CSIRO, Macquarie University and Swinburne University of Technology are well placed to support this growing need. Australian access to data from SKA, ESO or other facilities would be significantly enhanced with these kinds of systems.

**Analysis on Science Platforms:** The growth in the typical size of science-ready datasets has grown significantly over the last decade. While it's still possible to work on data on a local desktop, it's becoming increasingly difficult to move data around and to have the resources to analyse it efficiently. This is driving the move towards science platforms, which allow researchers to perform science analysis interactively via their web browser, on a secure computing facility co-located with the data.

There are good examples of science platforms internationally. For example, CANFAR specifically serves the astronomy community, providing specialised visualisation and analytics services via Jupyter notebooks and remote desktop access that rely on CANFAR's back-end HPC resources. Astronomy Commons serves the Zwicky Transient Facility within the ZTF Partnership and is built on Amazon Web Services, utilises Apache Spark for parallel data analytics, and JupyterHub as the web-accessible front-end.

Within Australia, the Australian Research Environment hosted by NCI provides similar services to CANFAR (e.g. JupyterLab, RStudio App, direct ssh access to NCI Gadi) but serves the range of numerate disciplines in Australia, including astronomy but also oceanography, climate modelling, etc… With an astronomy focus, Data Central is developing an Australian equivalent to CANFAR that should be ready in the second half of 2024. The SKA Regional Centres will be a global network of data centres which will provide users with SKA data access, user support, and a science platform for collaborative science. In Australia the Australian SKA Regional Centre (AusSRC) is already under development, with capacity and resources ramping up for Australian users to meet the SKA as it comes online.

While science platforms require significant initial investment and planning, they should enhance scientific productivity and outcomes. To remain competitive internationally, Australia requires one or more Science Platforms and the ability to link between them and with international counterparts.

**Mix'n'Match Astronomical Data:** The example of GW170817, the merger of two neutron stars, has shown that the ability to study astronomical objects across wavelengths and messengers (i.e. photons and gravitational waves) enables transformative and breakthrough science. Data fusion - bringing multiple datasets

together in one place to facilitate joint analysis as seamlessly as possible - is becoming feasible but requires further investment in both funding and infrastructure. This means multi-wavelength and multi-messenger datasets (including not only gravitational waves but also data from high energy astrophysics experiments encompassing cosmic rays and neutrinos), but also includes theoretical and astrophysical simulations datasets. Clearly this demand to compare distinct datasets of quite disparate origins emphasises the need for data that adheres to the [FAIR principles](#), which can be achieved by [IVOA standards](#) compliance; it requires domain expertise to place data on the same footing when they are brought together (e.g. [world coordinate systems](#), consistent units); and it can be most seamlessly achieved within a science platform model.

**The value of long-term preservation:**  A study of research data in Australia found that universities are holding approximately ~150 PB of on December 2021 (cf. [Research Data Culture Conversation - A Macro View of Retained Australian Academic Research Data](#)) , and with astronomy archives ([Data Central](#), [MWA](#), [CASDA](#), [SkyMapper](#) and [OzSTAR](#)) totalling >40 PB at the time, this means astronomy data alone make up ~25% of Australia's research data. As Table 3 in the International and Space Facilities white paper makes clear, the data volumes will continue to grow as we access new facilities.

The impact of archival data can be profound. A search of the [NASA ADS database](#) reveals that the number of publications per year which use archival telescope data has more than doubled over the past decade. The [US Astro2020 Decadal plan](#) noted that publications from [archival Hubble Space Telescope data](#) have outnumbered those by the original proposing teams, and other major facilities are seeing the same trends. There have been many noteworthy discoveries from archival data. The Australian example is the discovery of Fast Radio Bursts. The first FRB (the Lorimer burst, [Lorimer et al. 2007](#)) was only possible from re-analysis of archival Parkes data, six years after the original observations. This opened up a whole new field of astronomy. The community strongly supports prioritising long-term data storage, to ensure data resilience and security, as astronomical data has long-term scientific value.

**Recommendations:**
- **Standards**
  - Ensure that all Australian-hosted data are [IVOA standards compliant](#) and adheres to [FAIR Principles](#)
  - Ensure that data portals are designed to enhance accessibility and to improve user experience, for example, by adhering to international standards, [Web Content Accessibility Guidelines 2.1.](#)
- **Software**
  - Invest in Australia's data reduction pipeline capabilities - including scalability and automation - and improve funding for development and long-term support of pipelines.
  - Establish a mechanism for long-term astronomical software sustainability to ensure the impact and legacy of these codebases
- **Platforms**
  - Support key national research data platforms to provide long-term preservation and availability of data, via ongoing funding and regular capital refresh.
  - Invest in an Australian astronomy-focussed science platform, allowing interactive analysis of remote datasets that fuses together multi-wavelength, multi-messenger, and theory/simulations data.
- **Interoperability**
  - Ensure multi-wavelength and multi-messenger data centres are interoperable
  - Improve the discoverability and accessibility of astrophysical simulation data and its interoperability with observational data in astronomy

## 2. Community Perspectives on Software

**General Perspective:** The importance of software for the science that we do - how it's developed and maintained, and its ability to maximise the hardware it's deployed on - is widely recognised and broadly

understood. However, the typical astronomer does not receive formal training in modern software engineering and development practices, which is necessary to maximise software quality, scalability, and sustainability. While the subset of astronomy researchers with these skills has been growing over time, the community recognises the need for (1) general community-wide training in the essentials of software, computing, and data best practice, to enable researchers to make informed choices about software; and (2) support for software specialists with domain expertise.

**Specialised Astronomy Software Support:** A significant development since the last decadal plan is the emergence of the [Astronomy Data And Computing Services (ADACS)](). ADACS began in 2017 and is managed by [Astronomy Australia Limited (AAL)](). It provides software engineering and development services via its [merit allocation programme (MAP)](), as well as bespoke on-demand training and year-on-year internship opportunities. Via [ADACS MAP](), software specialists work with astronomers on specific projects, either building new software from the ground up or modernising and/or improving an existing code base; projects can last one or more semesters, with a semester's worth of work corresponding to approximately 3-4 months. The community regards [ADACS]() as a success. It has fulfilled the recommendation from the Mid-Term review to "provide funding for commensurate training and education in data science and code development".

**Research Software Engineering:** Over the last decade, there has been a growth of [Societies of Research Software Engineers (RSE)](), including in [Australia & New Zealand](), which provide support and resources for those engaged in RSE roles - "those in academia who combine expertise in programming with an intricate understanding of research. Although this combination of skills is extremely valuable, these people lack a formal place in the academic system" . The community notes that the research software engineering landscape is more developed overseas than in Australia. For example, [UK Research and Innovation funds grants]() to support the design, implementation, and development of software that is dependable, efficient and maintainable; the [UK Software Sustainability Initiative]() is dedicated to improving software in research; and [NASA's Transform to Open Science (TOPS)]() initiative is a 5-year programme to build the infrastructure to train scientists and researchers in software best practice.

**Future Considerations:** The advent of exascale computing, with its need for efficient, scalable algorithms and compliance with green computing, and the enormous data volumes and processing requirements of the [Square Kilometre Array]() and [LSST]() will place stringent demands on software performance and scalability. There has been a move towards software adhering to the [FAIR principles]() - that it be findable, accessible, inter-operable, and reusable - in the same way that data does; for it to adhere to [IVOA]() software/tools standards to support processing/analysis of multi-wavelength large data sets; and for it to be [open source](). The growing importance of AI/ML and the likely growth in quantum computing will also place demands on verifying reproducibility.

**Recommendations:**
- **Training**
    - Provide general community-wide training in the essentials of software, computing, and data best practice, to enable researchers to make informed choices about software.
- **Standards**
    - Ensure that all Australian-developed software is [IVOA]() standards compliant; adheres to the [software FAIR Principles](); and is [open source]().
- **Software Specialist Support**
    - Ensure the continued funding of [ADACS]()
    - Support research software engineering in general (via professional bodies such as [Research Software Engineers Society of Australia and New Zealand]())

## 3. Community Perspectives on Artificial Intelligence, Machine Learning (AI/ML)

**General Comments:** AI/ML has emerged as a significant research area in Australian astronomy over the last 10 years. In the previous decadal plan, this area was treated as a subset of data-intensive astronomy. There

have been several targeted hires in AI/ML for astronomy, including the formation of an [AL/ML team at CSIRO ATNF Science](), the [CSIRO Collaborative Intelligence Project with ASKAP](), and the [James Webb Australian Data Centre at Swinburne]().

The particular focus has been on using [deep learning]() and 'classical' networks such as [convolutional neural networks (CNNs)]() (e.g review by [4]) but there are more recent developments such as [diffusion models](), [attention networks]() (transformers) and [graph networks](). [Large Language Models]() (LLMs) have become advanced and their use is now widely spread, from enterprise apps and search engines to co-pilots for programming. Machine learning and artificial intelligence has the potential to transform science and have many use cases in astronomy. It will be an essential tool to maximise the science outcomes from large astronomical datasets.

The growth of machine learning applications requires astronomers to collaborate with computer scientists and engineers, and it requires a science platform infrastructure to fully exploit and make full use of the ML/AI techniques. The astronomy community may benefit from partnering with industry where complex ML frameworks/pipelines are already deployed. ML methods have wide applications in other scientific domains, in industry and have a large potential for economic return. Furthermore, data science and machine learning is an increasingly attractive career which the astronomy community can leverage to inspire students to undertake astronomy.

**Global Growth:** The emergence of AI/ML in Australian astronomy over the last decade mirrors that seen internationally. A search for papers with the keywords "machine learning" on the [NASA Abstract Data Service]() illustrates this. As Figure 1 indicates, the total number of refereed papers has grown by a factor of 7.7 between 2016 (132 papers) and 2023 (1018 papers), while the number with at least one Australia-based author has grown by a factor of 4.6 (10 to 46) over the same period.
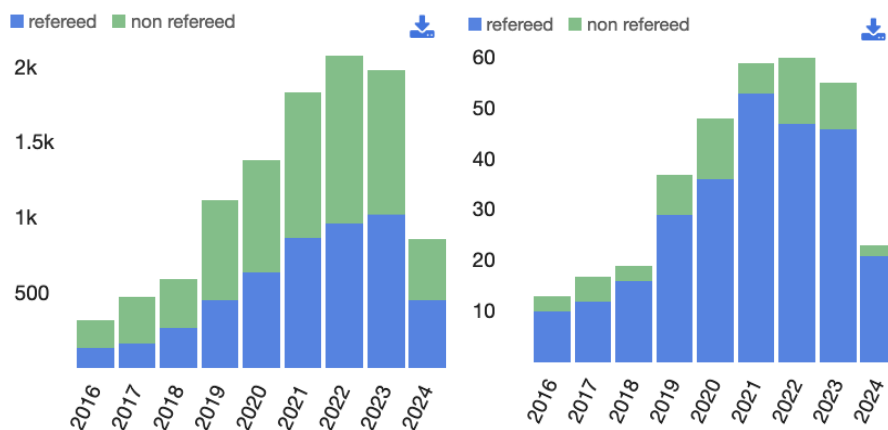


Figure 1: Left, all astronomy papers with machine learning keyword in 2016-24; right, with at least one author from Australia

The [United States Decadal Survey on Astronomy and Astrophysics (Astro2020)]() recognised the growing role of machine learning and data science, and mentioned the need for training at the graduate level. However, large language models have really only emerged in the last 4 to 5 years, after the decadal plan process, and ML arguably plays a bigger role now. There is significant investment in the US with the [National Science Foundation (NSF)]() and [Simons Foundation]() each committing [$US 20M to fund up to two ML/AI in astronomy institutes over 5 years]() (at $4M per year). [NASA]() has formed a working group to begin looking at [AI initiatives](). In Europe and the UK, there are several funded [European Research Council (ERC)]() and [UK Research and Innovation (UKRI)]() projects in astronomy which use machine learning, but nothing to the scale proposed by NSF and Simons Foundation.

In the last decade some investment and progress in this area has been made in Australia, with some targeted hires in individual institutes and universities for ML/AI projects. A meeting was held at the [2022 ASA Annual General Meeting](#), which initiated an informal "ML/AI in Astronomy" interest group and regular meetings were held but this group is now inactive. About 20% of respondents to the Decadal Plan community survey say they have made use of ML/AI, and 8% say it is critical for their research. About 30% of the community anticipate using ML in the future.

Australia is a relatively small country in the machine learning and artificial intelligence field compared to countries such as the United States, China, United Kingdom and Canada (see e.g. [see Australian Federal Government report on AI technologies](#)). However the Federal Government has recognised AI as a critical technology and invested in a [National Artificial Intelligence Centre](#). We can not compete with the leading countries in this field but our comparative advantage is in our large multi-wavelength and multi-messenger datasets.

**The next decade for ML/AI in astronomy:** In the next decade major observatories will come online and ML/AI techniques are required to fully exploit the large datasets from these facilities. For example,

- The [SKA](#) is currently under construction in Australia and South Africa, with science verification to begin in 2027 and science operations expected to start from 2029/2030. The [SKA](#) is expected to produce a torrent of data that can not be inspected by eye. Machine learning algorithms will be required for source classification and anomaly detection. Much work has already been done in this area for the classification of radio galaxy morphologies for [SKA](#) pathfinders (e.g. [Galvin et al. 2020](#), [Mostert et al. 2022](#), [Gupta et al. 2024](#)) and detection of pulsars ([Liu, Y et al. 2024](#), [Bhat. S. et al. 2023](#), [Lousto, C et al. 2022](#)).
- The [SKA](#) and its [Regional Centres](#) are expected to archive 100s of PB per year of science data. A potential high impact use case is data compression: machine learning models can be used to learn the features of a dataset without losing information. The climate change science community has similar big data challenges and is already exploring compression of multi-dimensional climate data with neural nets ([Liang, X et al. 2022](#), [Huang, L. and Hoefler T. 2022](#)).
- [Euclid](#) was launched in July 2023. It will conduct two major surveys, the [Euclid Wide Survey (EWS)](#) of about 15,000 square degrees and the [Euclid Deep Survey (EDS)](#) totalling about 50 square degrees. It is expected to obtain photometry of about 10 billion sources, shapes for 1.5 billion for weak lensing studies, and spectroscopic redshifts of 35 million galaxies for studies of clustering, evolution of large structures and the nature of dark matter. The complete survey will contain hundreds of thousands of images and several 10s of PBs of data, providing a rich dataset for machine learning algorithms. Un-supervised methods can group similar galaxies for the study of their physical properties. The millions of spectroscopic redshifts provides an excellent reference, or ground-truth, dataset for supervised learning of redshifts. In principle a neural network can learn the shapes of the galaxies in the data, which is one of the primary goals of [Euclid](#).
- The [Vera C. Rubin Observatory (VRO)](#) currently under construction on Chileon Cerro Pachón in Chile, will reach first light in 2025. It will undertake the [Legacy Survey of Space and Time](#) to survey the southern sky deeper and faster than any wide-field survey to date. It is expected to produce about 10 million alerts per night, enabling the discovery of an unprecedented large number of astrophysical transients and opening a new era of optical big data in astronomy. [FINK broker](#) (Moller et al. 2021) was one of only 7 brokers (software systems that will ingest, process, and serve astronomical alerts from the [VRO](#) and other surveys to the broader scientific community) chosen by the observatory in a competitive call process. software systems that will ingest, process, and serve astronomical alerts from the [VRO](#) to the broader scientific community. [FINK broker](#) science modules employ machine learning techniques to classify the transient source and detect anomalies. It is designed to adapt to new labels and be re-trained as the survey progresses.

The same is true in gravitational wave astronomy (e.g. [Cuoco et al. 2020](#)) and in theoretical astrophysics, where emulators (e.g. [Conceição et al. 2023](#)) and simulation based inference techniques (e.g. [List et al 2023](#)) are becoming commonplace.

Foundation models have transformed ML/AI and are powering today's commercial/industry generative models such as [ChatGPT](#). These large-scale models, trained on huge datasets, are capable of performing a wide range of tasks. They have the potential to accelerate discoveries by automating data analysis, identifying patterns, and making predictions. The large datasets coming in the next decade (mentioned above) have the potential to form the basis of transformational foundation models for astronomy, and work is already happening to develop astronomy foundation models using today's datasets (e..g [AstroCLIP](#): Parker et al. 2024). Foundation models embedded into the data archives themselves would lead to a revolutionary change in how users access and analyse data.

It's worth noting that the emergence of LLMs will impact how our community publishes papers. For example, there is growing interest in LLMs as a tool to generate citation recommendations (e.g. [Wu et al. 2024](#), [Iyer et al. 2024](#)). Major academic publishers now have [policies](#) in place around the use of LLMs to write text, and there are concerns around [plagiarism](#) and the use of [published work to train models](#). This is a rapidly evolving area and some careful consideration should be given to these issues.

Telescope operations is another area which will benefit from machine learning and artificial intelligence. Applications of ML to short and long term scheduling can improve telescope efficiency, and for example, may be particularly useful in following up the large number of [LSST](#) alerts. ML/AI can be used for predictive maintenance, especially in the era of distributed and modular telescopes such as the [SKA](#). Next-generation adaptive optics systems can benefit from the computational efficiencies of AI/ML methods for real time control.

Machine learning will benefit from science platforms with fast data access and good workflow management. Good data curation and adhering to [FAIR principles](#) will facilitate the use of astronomy data in ML applications. The bulk of the training and development of ML models is done on GPUs, and with industry ML applications driving the progress in GPU technology there will likely be further developments to GPU technology in the next decade (faster and more energy efficient GPU nodes, with more memory, for example). The growth of quantum computing may boost the performance of machine learning algorithms and lead to another revolution in ML/AI similar to the advent of GPUs, but it is unclear whether this technology will mature or progress enough in the period of this decadal plan.

**Recommendations:**
- **Training**
  - Provide training in the essentials of machine learning as a part of PhD training
- **Community of practice**
  - The Australian astronomical community should have an active formal community of practice on the topic of machine learning, to share knowledge and drive growth in ML uptake and trust
- **Funding of projects**
  - Large scale projects and funding proposals, e.g. ARC CoEs, should consider specifically funding ML/AI initiatives, where appropriate, for targeted outcomes, but any large significant investment in ML/AI is unlikely to complete with that of other countries

# 4. Community Perspectives on Computing

**General Comments:** For the purposes of this white paper, computing refers to high performance computing (HPC). HPC capability can be categorised into tiers, following the [PRACE](#) [convention](#)
1. Tier 0: international-scale, largest available for scientific research;
2. Tier 1: national-scale, beyond budget of an individual university or research institute;
3. Tier 2: individual university- or institute-scale.

Our focus is on HPC facilities that are open to the entire Australian astronomy research community. Currently Australian astronomers have access to two Tier 1 facilities ([Gadi at NCI](#); [Setonix at the Pawsey Supercomputing Centre](#)); two Tier 2 facilities ([OzStar](#) and [Ngarrgu Tindebeek (NT)](#) at [Swinburne University of Technology](#)); but no Tier 0 facility. Some universities provide access to tier 2 facilities and extra access to tier 1

facilities to their own researchers, but these are not community access. Similarly, the [SKA Observatory](#)'s Science Processing Centre in the [Pawsey Supercomputing Centre](#) will include the [Science Data Processor](#), a 135 Petaflop[1] machine, but this is dedicated to [SKA](#) data processing and is not community access.

Australian astronomy's HPC user base spans those who exclusively use tier 2 facilities to those who are among the largest users of Australia's Tier 1 facilities. Tier 1 facilities are accessed primarily through the [National Computational Merit Allocation Scheme (NCMAS)](#), which is held annually and awards computer time across all computationally intensive disciplines. [Astronomy Australia Limited (AAL)](#) buys additional time on [NCI Gadi](#) that it provides to the community via the [Astronomy Supercomputing Time Allocation Committee (ASTAC)](#) in two calls, one for large projects and one for smaller projects. The Tier 2 facilities [OzStar](#) and [NT](#) are accessible to all astronomers in Australia upon request, and astronomy usage comprises 50-70% of these facilities. Most time is available on a first-come, first-served basis, although [AAL](#) does provide a small amount (approximately 2 million CPU-hours per year) via a competitive grant program.

Large users - typically those who work on astrophysical simulations or radio astronomy processing - often combine awards of time from multiple sources (e.g., [NCMAS](#), [Australia Large Computing Grants](#) – ALCG, [AAL's ASTAC schemes](#)) to obtain awards larger than would be possible via a single scheme. A 2021 survey by AAL's Australian eResearch Advisory Committee found that approximately 50% of the total computing time used by Australian astronomers was on facilities outside Australia[2], which they accessed through a variety of schemes, including formal applications to schemes overseas that are open to non-residents; access through prior positions or joint appointments; and access via overseas collaborators.

**Open Access to HPC:** The total amount of computing time on tier 1 machines available to all numerate disciplines via [NCMAS](#) has increased in recent years, in particular in response to the availability of both phases of [Setonix at Pawsey](#). In the NCMAS round for 2024 allocations, the time available was,
- 150M core-hours on [NCI Gadi](#), the majority of this on CPU but a small part GPU
- 325M CPU core-hours on [Pawsey Setonix](#)
- 1.25 million GPU-hours on [Pawsey Setonix](#)

The oversubscription rate in this round was approximately a factor of 2. Astronomy applicants were awarded roughly 10% of computer time, totalling approximately 50 million CPU-hours and 100,000 GPU-hours; the largest individual awards to astronomy applicants were in the range of 10-15 million CPU-hours. An additional 11 million CPU-hours on [NCI Gadi](#) were available to astronomers via the [AAL ASTAC](#) rounds during the same period. The tier 2 machines provide of order
- 130M CPU hours per year (33M on [OzSTAR](#), 95M on [NT](#))
- 1M GPU hours per year (0.17M on [OzSTAR](#), 0.77M on [NT](#)),

of which time 50-70% is open access to the astronomy community - approximately 65M to 90M CPU hours and 0.5M to 0.7M GPU hours. When combined, the total compute resource openly accessible to the Australian astronomical community is roughly 125-150 million CPU-hours and 0.6 - 0.8 million GPU-hours per year. If we include the ALCG scheme, this increases the CPU component to 150-175 million CPU-hours.

The benchmark set in the last Decadal Plan was 30% of a top-100 machine. As of the time of writing (August 2024), [NCI Gadi](#) is ranked #103, [Pawsey Setonix](#) is ranked #28 on the [top-500](#) June 2024 list. If we consider their combined total compute capacity as equivalent to a top-100 machine, 30% is of order 1200 million CPU-hours and 3.6 million GPU-hours per year[3]. By this metric, the total resource available to the Australian HPC community is approximately 15% of the need anticipated in the last Decadal Plan (30% if we consider only [Pawsey Setonix](#) as the Australian machine within the top-100, and compare to the upper value of 175M

---

[1] 1 petaflop is a quadrillion, i.e. $10^{15}$, floating point operations per second.
[2] The primary reasons given for use of facilities outside Australia were that Australian facilities did not provide enough computer time, that machines outside Australia were significantly faster or otherwise more capable, and that the proposal burden (i.e. the amount of effort required to prepare a proposal) was significantly higher for access to time in Australia than for time outside Australia.
[3] Using 250k CPUs and 640 GPUs for [NCI Gadi](#), and 200k CPUs and 750 GPUs for [Pawsey Setonix.](#)

CPU-hours), or about 5-10% of a top-100 machine. By this measure, the recommendation of the previous Decadal Plan has not been met.

Compared to peer communities overseas - those with similar GDPs and research community sizes, such as Canada - the total compute capacity per GDP is somewhat below the average (see Table 2 in appendix), but the amount of **open access** compute time is well below the average. This has led to many large Australian HPC users to access more compute time via overseas collaborators or other connections to make their work possible. This lack of access is partly a result of a smaller HPC infrastructure base, but a more important contributor is HPC access policies, which in Australia reserve the great majority of HPC time for dedicated projects, leaving only a small fraction available in the open merit-based programs from which most astronomical HPC applications draw. The comparison with Canada is particularly telling: its total HPC capacity per unit GDP is only 70% of Australia's, but it provides 2-3 times as much open access time. It can be argued that the current lack of access strongly inhibits the development of the Australian simulation community.

Australia's tier 2 facilities that target astronomical computing - OzSTAR and NT - are very well regarded by the community - they are well managed, they work well, they provide a solid base of HPC for smaller users, and there is a desire to see their continued funding with regular capital refresh in the coming decade.

**Radio astronomy's HPC requirements:** HPC is a critical component of the operating model of modern radio telescopes, providing the essential processing backend that performs correlations and beamforming. The Pawsey Supercomputing Centre plays a central role, reflecting its importance role as a member of the SKA's Science Data Processor (SDP) Consortium, and the investment in Setonix (#28 in the top-500's June 2024 list) has been welcomed by the Australian radio astronomy community. Since the previous decadal plan, the goal of moving processing from a user- to an observatory-operated processing model has been achieved. ASKAP operational processing has moved from Pawsey Galaxy to Setonix and is now close to a sustainable automated workflow; there are 180 dedicated nodes providing approximately 170M hours of compute per year. MWA operational processing takes place on Pawsey Garrawarla. Reaching full survey operations has taken longer than anticipated, in part because of insufficient compute during the initial stages, in part because of delays in the deployment of Setonix and its integration into the operational workflow. The ASKAP operational workflow is producing science products but has yet to reach full capacity. Ongoing operational support will be essential as ASKAP and MWA surveys run at full capacity - both to process the raw data but also to facilitate serving of Petabyte-scale datasets - and a stable and reliable system is critical for the real-time, high throughput compute demanded by ASKAP, the MWA, and eventually SKA-Low. As SKA-Low comes online, the planned 135 Petaflop Science Data Processor should provide 27B CPU hours per year.

Science user access for non-operations compute is maintained through merit allocation access to Setonix, which has worked well in providing additional compute resources for the ASKAP community. As data volumes grow over the course of the next 5 years, especially as the ASKAP surveys achieve full operations performance, the processing requirements will become more advanced and produce more data that needs analysis, and there is a case for this processing to become part of the operational workflow. This also places unique demands on operations-critical compute platforms for high-data-rate telescopes, because of the sustained data throughput and the emphasis this places on storage, filesystems and isolation from general users to ensure reliability; this should be accounted for when undertaking a capital refresh.

Although a capital refresh is expected during the period of the next Decadal Plan, it is unlikely that the increase in computing power between Setonix and its successor will match that between Galaxy and Setonix. In addition, demand from new science domains means that the successor will have a larger user base, which means that compute resources for radio astronomy may not grow. This provides strong motivation for software optimisation and refinement, and investment in software - and the specialists that develop and deploy it - that maximises its performance on the resources it has access to. There is as yet insufficient investment in preparing existing software packages and analysis workflows to run at scale by making use of (multi-)GPUs or even more basic HPC technologies like MPI, and insufficient investment in instrument-dedicated hardware -

although ASKAP has bucked that trend and invested significant effort in developing YandaSoft that can run at scale on HPC systems by using MPI. It's crucial in the coming Decadal Plan period to invest in the platform and software engineers that will develop portable, scalable codes and workflows to make best use of the increasingly heterogeneous systems that will underpin future HPC facilities. This will require a more collaborative approach between radio astronomy, HPC centres, and researchers in the field of HPC, data science and AI.

AusSRC will play a critical role in providing access to the necessary HPC resources for doing science. Although it will directly support SKA observations and analysis, the transition to SKA operations presents an opportunity to provide additional resources to help radio astronomy research more widely, in addition to resources available at Pawsey and NCI. It will be important that these resources keep pace with the growing size of the SKA arrays (with an emphasis on SKA-Low, but also SKA-Mid). They will also need to be appropriately identified and put in place in a timely fashion to avoid lost opportunities.

**Key Recommendations:**
- **Open-access HPC**
  - The astronomical community should continue to support OzSTAR and NT  and the programs initiated by AAL to provide computing time on NCI. AAL should consider extending this access program to the Pawsey Supercomputing Centre, and should investigate other options including commercial providers and overseas facilities.
  - Australia should move more of its computing resources into an open, merit-based allocation scheme, with the goal of reaching approximately 30-40% open access – in line with peer countries – over the decade. This will benefit not just Australian astronomy, but science and engineering in Australia more generally.
- **Large-scale HPC programmes**
  - There should be enough HPC time available for it to be possible to carry out large programs comparable to those that are possible in peer nations - a target is 30% of a top-100 machine is recommended.
- **Specialist software support**
  - The astronomy community should support a specialised software service to help researchers migrate from CPU- to GPU-based architectures, which increasingly dominate the top end of HPC hardware, and to maximise software and workflow performance on heterogeneous architectures in current and future HPC facilities.
- **Supporting the SKA and its precursors**
  - The continued support of ASKAP and MWA by stable and reliable HPC systems is critical for the real-time, high throughput compute required by the processing and analysis workflows to support surveys in the coming decade.
  - Support AusSRC to provide resources to support radio astronomy research, in addition to resources available at Pawsey and NCI, ensuring that the resources keep pace with the growing size of the SKA arrays and are in place in a timely fashion to avoid lost opportunities.

# 5. Community Perspectives on Skills & Culture

**General Comments:** There is a widespread recognition that data and computing are key pillars of astronomical research, which is important because a growing proportion of the Australian astronomical community requires a high performance computing and/or data infrastructure to carry out their work. The complexity and scale of this infrastructure means that dedicated specialists - in software, data, and computing platforms - with domain expertise (e.g. in radio astronomy or numerical simulations) play a crucial role in facilitating research. This has underpinned the growing trend in astronomical software and data away from "the individual researcher with laptop" towards a larger-scale team-based approach. The community does well at providing training opportunities, especially for PhD students and early career researchers, to acquire the necessary skills to work with this infrastructure. Sustained and close collaboration between researchers and

infrastructure specialists will be essential in the coming decade, which is in line with the Australian Federal Government's National Digital Research Infrastructure Strategy. This will also require building career pathways for these infrastructure specialists.

**Computing & Data Competency**: The general researcher should have,
- Basic Linux/unix skills and an aptitude in an open source scripting language, such as python or R;
- Understanding of the basics principles of computer programming and how to work with different data structures, and of code optimisation and the factors that affect code performance;
- Familiarity with the basic principles and practice of source and version control, using tools like git and repositories such as GitHub or GitLab; and
- Ability to use notebooks (e.g. Jupyter) to work with remote datasets.

These core skills should form the basis of a nationally coordinated training and accreditation scheme for data & computing competency, incorporated into PhDs.

The researcher who requires more advanced skills should have a mastery of,
- A compiled language (e.g. C/C++, Rust, Fortran, Java) and a build system such as cmake or meson;
- Knowledge of parallel programming models (MPI, OpenMP, CUDA/HIP) and IO technologies (HDF5, ADIOS2, CAPIO) critical for effective use of HPC/D, as well as technologies such as cuPy, mpi4py and Numba that can provide much of the performance of CUDA and MPI but can be developed within a scripting language framework.
- Queue based batch systems (e.g Slurm, PBS) that are commonly used in HPC/D environments;
- Performance monitoring and profiling;
- Basic principles of software engineering such as code commenting and documentation, unit tests, git merging branches and code rebasing; and
- Understanding of system architecture and its impact on code performance.

It's worth noting that as projects reach a threshold of complexity, there should be explicit collaboration between the domain specialist researchers (with the requisite knowledge in data and computing) and the specialist data, software, and platform engineers who can build the codebases and workflows adhering to best practice.

**New Ways of Working:** The growth in the rate and volume of data generation makes it impractical to copy large datasets between remote machines and local machines for processing and analysis. This should prompt a shift towards *bringing code to the data* - where data analytics hubs are attached to computing or data centres and allow users to access their data via notebook-style environments. A good example of this is provided by the Australian Research Environment (ARE), which is hosted by NCI. ARE is a web-based graphical interface that can access the NCI Gadi supercomputer and NCI's data collections, offering applications such as virtual desktops and JupyterLab. The IVOA is moving towards running jobs remotely, offering a set of tools and standards; given how the IVOA is widely used by all significant surveys, there is an opportunity for the Australian community to take advantage of this by ensuring that homegrown software and datasets are IVOA-compliant.

Software and data requirements are now sufficiently large and complex that **communities of practice** are necessary - groups of researchers within specific sub-disciplines that collectively manage their shared software and data needs, working with teams of dedicated infrastructure specialists. A model for this already exists in the form of ADACS and the Australian SKA Regional Centre, in which dedicated teams of software developers are central to strategic software initiatives, skilled in project management and software design best practice. This provides a template for projects that need to maximally utilise HPC/D facilities - where researchers work with properly resourced, dedicated teams with specialist expertise.

Collaborative code development is well established within astronomy, with tools such as git and repositories such as GitHub and GitLab playing a central role. It has practical benefits - e.g. it allows for branch-controlled workflows with version control, automated testing, and backlog management (with e.g. Jira, GitLab) - but it also plays an important role in scientific reproducibility.

Given the fundamental importance of software in modern research, more thought needs to be given to software security - how do we ensure robust, resilient, sustainable software with long term support? Arguably this is a weakness - anecdotally not enough of the software that underpins our research activities is secure by this metric. Many packages are maintained by small teams - and subject to single point failures when e.g. critical skills are lost when an early career researcher moves on. NCRIS and NDRI have the opportunity to better support widely-used mission critical software; investment in software and data support should be a condition for funding of new telescopes and instruments.

Over the next decade, we also expect to see a significant rise in the prevalence of mixed/extended reality (MR/XR) technology that is directly applicable to scientific contexts. To date, we have seen the application of virtual and augmented reality (VR/AR) mainly in specific localised aspects of astronomical data/computing such as multi-dimensional data visualisation, outreach/science communication/education and in some cases of collaboration/meetings (e.g. Fluke & Barnes 2018, Hiramatsu et al. 2021, Jarrett et al 2021, Impey & Danehy 2022, Kersting et al. 2022, Endeavour 2023, Moss et al. 2023). The rapidly advancing technology along with increasing affordability will likely result in a widespread application of XR across the domain of astronomy (also explored as part of the WG 2.2 report in the specific context of national and university facilities), offering new opportunities for the next generation of astronomers to explore, analyse and collaborate on their data.

Finally, as noted in the community perspectives on AI/ML, LLMs have evolved rapidly in the last couple of years and will impact how our community publishes papers, from the benign (generate citation recommendations) to the malign (potential for plagiarism). This is a rapidly evolving area and some careful consideration should be given to these issues.

**Stable Long-Term Career Pathways:** A consistent picture has emerged from community feedback of the critical role of data, software, and platform specialists in providing the foundations for the Australian astronomical community's research programme. At present researchers who have developed expertise in these areas and who wish to specialise in them face a precarious career if they remain in astronomy research, given the limited opportunities for career progression. There is a strong community desire for there to be credible pathways for these specialists to build stable and long-term careers in astronomy; without them, these specialists will bring their highly prized expertise elsewhere.

There are examples of universities that offer opportunities for researchers who focus on e.g. research software development to pursue a career within the university environment but distinct from the traditional academic and professional staff roles (see, for example, the Melbourne Data Analytics Platform, MDAP). Similarly, there are dedicated astronomy specialists attached to the likes of ASVO, ADACS, Pawsey Supercomputing Centre, and AusSRC. However, the scale of the data and computing requirements in the coming decade mean that the Australian astronomy research community will need more specialists, and for these careers to offer the stability and opportunities for progression that will provide job security for the researchers and a skilled cohort for the research community. The importance of this is recognised in the emerging National Digital Research Infrastructure Strategy, and there is an opportunity for astronomy to address this need in the coming Decadal Plan period.

**Recommendations:**
- **Stable Careers**
  - Stable and predictable pathways for career development and long-term security are essential for the data, software, and platform specialists that underpin the astronomical community's data and computing needs.
- **Training**
  - Establish a nationally coordinated training and accreditation scheme for data & computing competency, incorporated into PhDs.
- **Blended research teams**

- Complex data and software projects should be explicit collaborations between domain specialist researchers and specialist data, software, and platform engineers who can build the codebases and workflows adhering to best practice.
- **Bring Code to the Data**
  - Invest in infrastructure to bring code to the data - where data analytics hubs are attached to computing or data centres and allow users to access their data via notebook-style environments.
- **Communities of Practice**
  - Establish communities of practice - groups of researchers within specific sub-disciplines that collectively manage their shared software and data needs, working with teams of dedicated infrastructure specialists.
- **Software Security**
  - Invest in widely-used mission critical software with explicit funding to ensure software security; this should be a condition for funding of new telescopes and instruments
- **Technology adoption**
  - Effective adoption of advancing novel technologies such as XR/VR/AR to enable new and innovative pathways to scientific outcomes alongside efficient distributed collaboration methods

## Notes on Technology

Interestingly, feedback on new technologies focussed on near term challenges.

- The emergence of ARM/RISC-V CPUs, which are being pushed by Nvidia as an alternative to x86, should bring enhanced performance but will require non-negligible effort to ensure that native (i.e. non-Python) code works correctly.
- The chiplet architecture of current and next generation AMD CPUs and upcoming Intel CPUs may mean that shared memory codes may not scale as well as expected.
- The performance of high-throughput filesystems like Lustre/GPFS are significantly adversely impacted by workflows that produce lots of small files, which is characteristic of radio astronomy.
- The power requirements of GPUs are beginning to exceed 800W per device, and similarly for CPUs. This is driving an increase in compute density, which will require modifications of both software and workflows.

There are opportunities.

- The emergence of a diversity of accelerators driven by AI will provide an opportunity for vastly increasing the scope of simulations and the variety of simulations that can be done.
- Data processing should also benefit from AI. Astronomy has the unique need to develop ML techniques that require minimal training data with uncertain ground truths, which is not the standard application of AIs. This unique need presents a unique opportunity.
- Quantum computing offers the potential for a low-energy, low carbon-footprint way of doing machine learning - but this is unlikely to happen on a 10 year timescale.
- Rise in distributed and remote work on global scales, to be increasingly facilitated by improved technologies for more effective collaboration across distance (e.g. mixed/extended reality)

Green computing has come to the fore in recent years as computing's contribution to carbon emissions via its energy consumption by high performance computing and data centres has been recognised. IBM reports that the information and communication technology sector is responsible for between 1.8% and 3.9% of global greenhouse gas emissions, while data centres account for 3% of annual total energy consumption — an increase of 100% in the last decade. Stevens et al (2020) estimated that "Australian astronomers' total greenhouse gas emwions from their … supercomputer usage" exceeded 15 kilotons of $CO_2$ per year, a significant fraction of the total $CO_2$ budget. Pawsey's Setonix is a green supercomputer, ranked #4 in the world in 2022 - made possible by a mixture of GPUs, innovative direct liquid cooling of all system components, geothermal cooling, and solar cells. The most direct impact on astronomy will be the need to move to GPU-enabled codes, to take advantage of the latest generation of supercomputers, and the need for

professionalised software development and data management to ensure these the facilities our community has access to will be used as efficiently as possible.

# Appendix

**Approach:** We identified 5 topics - high performance computing (HPC), data, software, AI & machine learning, and culture - and asked each group to consider these core questions, modifying them as appropriate for the topic.

1. What progress has been made in this area in the Australian community since the last decadal plan? Please refer to that plan and the status update in the mid-term review.
2. Has the level of progress matched the specific requirements of the Australian community?
3. How does this progress compare to that seen in comparable communities overseas?
4. Looking forward to 2026 to 2035, what level of progress is required to satisfy the community's requirements?
5. Are there particular requirements that need to be prioritised?
6. Are there particular opportunities that need to be accounted for when planning for 2026 to 2035?
7. Are there particular risks that need to be accounted for when planning for 2026 to 2035?
8. How does this compare to planning by comparable communities overseas?

**Table 1: Sub-Working Group Leads**

| Topic | Lead |
|---|---|
| High Performance Computing | Prof Mark Krumholz (ANU) |
| Data | Dr Simon O'Toole (AAO/Macquarie) |
| Software | Dr JC Guzman (SKAO) |
| AI & Machine Learning | Dr Minh Huynh (CSIRO) |
| Skills and Working Culture | Prof Chris Power (UWA) |

**Table 2: HPC Capacity in Peer Countries** (Credit: Mark Krumholz)

| Country | # systems | Total capacity (TFLOPS) | GDP (trillions of USD) | Capacity/GDP | Capacity/GDP relative to Australia |
|---|---|---|---|---|---|
| Australia | 6 | 48,784 | 1.7 | 28,696 | 1 |
| Canada | 10 | 41,208 | 2.1 | 19,622 | 0.68 |
| Germany | 36 | 256,266 | 4.1 | 62,503 | 2.18 |
| Switzerland | 3 | 26,667 | 0.82 | 32,540 | 1.13 |
| UK | 15 | 81,707 | 3.1 | 26,357 | 0.92 |
| USA | 160 | 3,725,850 | 25.4 | 146,687 | 5.11 |

Additional statistics on HPC in the European region can be found at [EuroCC-Access](EuroCC-Access).

**Table 3: Panel Members**

| Member | Participation |
| --- | --- |
| Arash Bahramian (ICRAR/Curtin) | WG1.3 liaison; Data |
| Krzysztof Bolejko (UTAS) | Culture and Skills |
| Sven Buder (ANU) | Data; Software |
| Lachlan Campbell (CSIRO Pawsey) | HPC; Data |
| Liz Davies (Data Central) | Data; Software |
| Celine D'Ogreville (ANU) | Data |
| Pascal Elahi (CSIRO Pawsey) | HPC for Radio Astronomy; Software |
| Christoph Federrath (ANU) | HPC |
| JC Guzman (SKAO) | Software WG lead |
| Roger Haynes (ANU) | Data |
| Alexander Heger (Monash) | HPC |
| Aidan Hotan (CSIRO Astronomy and Space Science) | WG2.2, 2.4 liaison; AI/ML; HPC |
| Jarrod Hurley (Swinburne) | HPC |
| Minh Huynh (CSIRO Astronomy and Space Science) | AI/ML WG lead; WG2.2 liaison; Data |
| Mike Kriele (ICRAR/UWA) | HPC; Data |
| Mark Krumholz (ANU) | HPC lead |
| Nuria Lorente (Data Central) | Software |
| Ridhima Nunhokee (ICRAR/Curtin) | WG1.1 liaison; Software |
| Sarah Martell (UNSW) | WG1.2 liaison |
| Richard McDermid (Macquarie) | Data |
| Vanessa Moss (CSIRO Astronomy and Space Science) | Skills and Culture |
| Simon O'Toole (Data Central) | Data WG lead; Skills and Culture; Software |
| Greg Poole (ADACS & Swinburne) | Skills and Culture; Software |
| Chris Power (ICRAR/UWA) | Skills and Culture lead |
| Aaron Robotham (ICRAR/UWA) | Data; Software |
| Manodeep Sinha (Sorsery Consulting) | WG3.1 liaison; Software |
| Marcin Sokolowski (ICRAR/Curtin) | Software |

| Steven Tingay (ICRAR/Curtin) | HPC and Data for Radio Astronomy |
|---|---|
| Emily Wisnioski (ANU) | Data |