

The Alan Turing Institute



Data Study Group: Challenge Owner Pack

Foreword

Thank you for joining The Alan Turing Institute, in its journey to make great leaps in data science and artificial intelligence in order to change the world for the better. By being part of the Data Study Group, you are, in a small way, helping us achieve our mission.

The Data Study Group is focused around the participants that apply to us, to work on your problems. These are the future leaders in the field, and we believe the more experience and opportunity they receive to experiment with real-world data, the greater they will become at solving the challenges of tomorrow.

Dr Sebastian Vollmer,
Data Study Groups Director

Contents

Introduction to the Data Study Group	4
Joining the pipeline	5
DSG schedule	6
What makes a good challenge?	7
Challenge descriptions and presentation	8
The short challenge description	8
The long challenge description	9
Top tips for challenge descriptions	9
Ethics and project sensitivities	10
The Monday challenge presentations	10
Challenge Owner deliverables	12
Challenge cycle timeline	13
Data security	14
DSG challenge data classification process	14
Loading data into the environment	15
The role of the Principal Investigator (PI)	16
Your first meeting with the PI	16
The final report	17
Attending DSG week	18
Get your hands dirty	18
Additional guests	18
Code of Conduct	19
Follow-on collaboration	19
Data Study Group contacts	20
FAQs	21
Annex 1: Example long challenge description	22
Annex 2: Sensitivity tiers and assessment	24
Turing Data Safe Haven tiers and flow	24
How to assess the tier of the projects	26

Introduction to the Data Study Group

Data Study Groups (DSGs) are intensive five-day collaborative hackathons hosted by the Turing. They bring together organisations from industry, government, and the third sector with talented multi-disciplinary researchers. Organisations act as DSG 'Challenge Owners' (COs) who put forth real-world problems to be tackled by small groups of gifted and carefully selected researchers.

During the DSG week, researchers brainstorm and engineer data science solutions. They present their work on the final day of the DSG and produce a report that will be returned to the CO and subsequently published on the Turing website.

The DSG is quite a large undertaking, not just for the participants, but also for the organisations that want to get involved. From the beginning of the engagement with us to the event itself, we estimate it takes about 3-6 months of preparation. COs must be fully engaged when shaping their challenges and take primary responsibility for providing adequate data.

Although the DSG is a big commitment, we offer constant support to COs at every stage of the DSG process. All business contracts will be handled by our business development team. The challenge Principal Investigator (PI) will work with a member of your team on all things scientific, data or challenge related.

We also have a dedicated IT team who will create and maintain individual environments for you to store and protect your data for the event.

It is important to note that participation in the DSG is not consulting. While we will work with you to design the challenge to meet your needs, we do so with the consideration of how this research will benefit others. The challenge posed is a framework for the participants who attend the DSG to explore your data in any way they see fit.

It is this freedom, coupled with no expectation to return positive results that attract them to get involved, and engenders possible novel solutions or directions for further research. The participants are not paid for this work.*



*Participants do receive a small travel allowance, accommodation if not living close to the event, and sustenance for the week.

Joining the pipeline

Before the challenge can progress to the next stage, all COs are invited to attend an initial challenge assessment. During this phase, the challenge is discussed with one of our data scientists to explore the feasibility of the proposal. Not all challenges are suitable for the DSG, and there may be other, more appropriate ways to engage with the Turing.

The assessment takes the form of an informal one-hour meeting between an experienced Turing Research Fellow and CO representatives. It is highly advised that CO representatives are able to talk technically about the challenge and the proposed datasets.

The challenge assessment consists of the following considerations:

1. Turning the initial challenge description into a well-defined scientific question
2. Aligning with the Institute's aims and themes
3. Framing as a three-day exploration and coding activity
4. Optimising the outcome of the DSG itself (e.g. report, impact) as well as the potential for a follow up project

So long as both parties are happy to proceed, the necessary contracts are drawn up, including the Collaboration Agreement. Once these have been signed, we recruit a PI for the challenge, and the real work begins.

The CO then works closely with the PI over the subsequent weeks to define and shape the challenge. During this period, the CO is also required to provide the dataset to be used during the DSG. The PI can advise on the quality and readiness of the data provided; however, it is down to the CO to process and prepare.

Due to the nature of the event, data needs to be as clean as possible, suitable for the challenge, and anonymised where there are sensitivities. More details on this are below.



DSG schedule

The DSG week is five days and is divided into three sections:

1. The Monday morning sales pitches
2. The work
3. The Friday group presentations

On Monday morning participants will hear from all COs – up to six per event – about their challenges and why they should get involved in their project. Following lunch, participants choose their groups and the work begins. Usually on Monday evening, we will take all participants (and COs if they want to join) to an offsite social event.

From Tuesday to Thursday the participants work hard in running experiments on your data, brainstorming avenues of research and writing up what they have done and how well it worked for the report and final presentation.

On Friday, the individual groups present back to the whole DSG, including all the COs. Here you will get a brief overview of what happened and what was (or wasn't) discovered.

Post event, we will take some time to prepare and finalise the report. The report will be published on the Turing website after about four months. You can find out more on reports and results further into this document. Publishing the report is important as we want the participants to be able to reference their involvement in the event.

As the national institute for data science and AI, it is also one of the requirements of the Institute to share the research generated here. You will be given the opportunity to review the report prior to publication to ensure that confidential information has not been disclosed.



What makes a good challenge?

A good challenge uses the strengths of the DSG scheme, in providing participants and the CO with an enjoyable and informative experience, as well as creating ample opportunities for impactful follow-up.

To ensure this, the PI and CO must work together to shape the challenge into something that is suited to the five-day setting of the DSG. Concretely, challenges should:

- Be realistic to explore within one day of brainstorming and three days of data science work.
- Be realistic to address with the data provided.
- Not have undue risks from data sharing, e.g. ethics, technical restrictions, privacy constraints, or data quality.
- Be well-specified enough to give participants a good start with low-hanging fruit, leading into more exploratory or less well-defined questions that may be more difficult.
- Focus on analytics and AI, rather than on rote tasks such as data munging, data curation, or data scraping.
- Be appealing to participants; with real-world impact, the potential to turn into a long-term project, or the right level of intellectual data-scientific challenge.
- Be likely to lead into impactful medium-term or long-term projects with Turing partners and participants, that can be kick-started by a DSG proof-of-concept or exploration.

The optimal trajectory of a challenge is as follows:

- Participants with the right skills read the long description of the challenge, and self-assign to the challenge they are interested in and feel they can contribute towards.
- During the week, the challenge team of CO representatives, PI, and participants produces proof-of-concept solutions for the low-hanging fruit challenges and brainstorms a series of approaches for the wider context.
- A report (and occasionally code) is produced which will be shared with the CO, as well as published on the Turing website.

In the context of wider research engagement:

- In the short-term, a follow-on project group forms around the suggested directions from the DSG. Longer-term and larger-scale project planning is informed by this seed research.
- Results of the follow-on project get published in major scientific venues, and/or lead to disruptive innovations which in turn inform further collaborative research projects, embedded in a long-term partnership network.

Challenge descriptions and presentation

In collaboration with the PI, we will need a challenge title and two ready-to-use descriptive pieces of text:

- Title of the challenge
- A one-paragraph short description of your challenge. This will be publicised when we open for applications
- An approximately two-page description of the challenge, to be used in the delegate pack for participants, and for the Turing's internal ethics approval process

In addition, the CO will be asked to give a 15-minute presentation (including five minutes to answer questions) of their challenge to the participants on the first day of the event.

The original challenge proposal form that was submitted at the very beginning of the DSG engagement can be used as the basis for these descriptions.

The subsequent sections provide high-level guidance for drafting the two challenge descriptions and the Monday challenge presentation slides.

The short challenge description

The purpose of the short descriptions is to attract participants to the data study group event, through various communication channels including the Turing website, e-mail marketing, social media etc.'

As such, it should be short, but also informative and address:

- What is the technical goal and scientific purpose?
- What data is being used, and who is providing it?
- What are the potential positive outcomes, i.e. research impact or business/societal benefit?
- What concretely is unique, special, once-in-a-lifetime, about the challenge?

It is important to note that the target audience is not the general public, but competent data scientists who may be interested in participating. Therefore, use of common technical data scientific terms (e.g. "supervised prediction", "randomised trial") is fine, and recommended where it helps create clarity.

The short description should highlight the challenge questions and associated opportunities – rather than, for example merely stating the general area, or stating how unique it is without saying why.

Exemplar short description:

Roche - 'Personalised lung cancer treatment modelling using electronic health records and genomics'

"Roche, Foundation Medicine and Flatiron are providing a recently collected, systematic and representative dataset comprising tens of thousands of US lung cancer patients' electronic health records, including detailed omics. Participants are invited to investigate whether modern data science and AI can help predict individual responses to different treatments, and how (or whether) these predictions can be leveraged for therapy recommendations."

The long challenge description

The long description will be provided to participants, for all challenges, 2-3 weeks in advance of the Data Study Group. Besides the Monday presentations, it is the main source that informs participants towards their decision to self-assign to one of the challenge groups within the week.

The second purpose is to provide full challenge documentation for the mandatory review by the Turing's Ethics Advisory Board, to be included in the project's ethics review submission to be carried out by the PI.

The two-page summary can have any structure. The content template to the right may be taken as a guide.

It should be written by the Challenge Owner, with the PI offering data scientific support.

Please see Annex 1 for exemplar long descriptions.

Overview section

- Short description of the challenge – a few sentences in length.
- Real-world context, overview/explanation to a general academic audience - 1/2 paragraphs.
- Who is the Challenge Owner? Possibly include a very short bio of the PI and company representatives - a few sentences in length.

Detail section

- Description of the data - a few sentences in length.
- Detailed description of the challenge, optimally with a very concrete low-hanging fruit sub-challenge to get people interested and stretch challenges - a paragraph and/or a couple of bullet points.
- High-level overview of potentially useful methods and/or approaches tried internally. Keywords to attract the “right” skills.
- Paragraph on potential follow-up activities, conclude with something inviting.

Top tips for challenge descriptions

- Start with the long description and once agreed, condense into the short description. The initial ethics approval will need to take place before the challenge is advertised for participant recruitment. Refinement to the long description can take place after the submission of the short description.
- Vague references or marketing language should be avoided entirely in both short and long descriptions.
- Language that is too prescriptive about the approach (“participants will apply [specific method X] in order to”), as opposed to descriptive of the challenge questions, should also be avoided.
- Remember, the challenge is a framework for participants to explore the data.
- Emphasise the wider context – what is the societal impact? How will the exploration of this data and challenge benefit more than just your organisation?

Ethics and project sensitivities

All research undertaken by the Institute is required to undergo an ethics review. The PI will collate the required information and submit it, but will need information from the CO. Ethical considerations when shaping the challenge are:

- Is the data sensitive, e.g. personally sensitive and/or subject to GDPR?
- If there is personal data, what is the process by which informed consent has been given?
- What are the risks of the project, posing potential harm to individuals or society?
- What measures of mitigation have been undertaken with respect to these risks?

While the PI submits the (Turing internal) approval request form, much of the information will require input that only the CO can give. Incorrect or inaccurate information about the challenge and data can delay the project and in extreme cases risk the viability of the whole engagement.

It is the CO's responsibility to have the correct authorisation to use the data and to know the intricacies of the data they are providing.

The CO should also prepare to assess the "sensitivity tier" and security arrangements for the hosting of their data and challenge. More on this in the Data Security section below.

The Monday challenge presentations

The Data Study Group event starts with the Challenge Owners' presentations of the challenges, to inform participants' self-assignment later on in the afternoon.

You will have 15 minutes to present and answer questions on your challenge (we advise 10 minutes for the presentation and five minutes for Q&A). There will be further opportunities for participants to ask you questions over the lunch break.

Participants will usually have read the detailed challenge descriptions, sent out the week before. It is nevertheless a good idea to make the presentation self-contained and explanatory to participants who are yet to decide.

The presentation should be aimed at a technical audience who has a relatively broad data scientific expertise but not necessarily in the domain. Therefore, make sure to explain basic domain terminology, but assume that the audience will be able to follow longer chains of reasoning that are necessary to

provide the technical explanation of the domain questions. The audience will expect and appreciate this approach.

In particular, please avoid giving the type of presentation that is exclusively targeting the general public, or an audience which is non-technical.

Please be aware that some participants are attracted by being able to work on a dataset that they wouldn't normally have access to and that will be locked-down - others are very committed to open science and may not be interested in closed challenges. Being transparent about these elements means you get the right participants for your specific challenge.

Again, also highlight how this challenge can benefit others outside of your organisation. What are the societal implications of making progress on your challenge? Showing the wider context of why your challenge is important will also draw participants to select your group.

The Monday challenge presentations (cont.)

Below is a suggested template structure, which of course may be deviated from:

- Introduce your organisation and the team. Introduce yourself and delegates for the week.
- A high-level overview of the domain application problem, and the wider context of the problem. Note that most participants will not be experts in the domain, but they will be experts in varied aspects of data science.
- A detailed, technical overview of the challenge:
 - The data – what data is being used, and how was it obtained? High-level summary of data schema and tables. What are the samples, how many, what are the variables?
 - Any data sensitivities, restrictions, conditions arising from the ethics review.
 - The challenge questions – what are the main questions? How does this tie into the general context? How does this translate into data scientific questions?
 - Avoid being prescriptive of approaches but be concrete in terms of goals.
 - Low-hanging fruit and stretch goals, leading into potential follow-on.
- Things you have tried before to address the challenge - things that worked or did not.
- What the future holds if this challenge (or greater problem it sits within) is solved, potential follow-on opportunities – concrete is better here, as future projects are a strong motivator.

Important note: we have found from experience that specifying what skills you might need for your challenge may discourage participants from other technical backgrounds from selecting your challenge. It is the multi-disciplinary aspect of the Data Study Group which can present you with some very novel approaches or solutions (at the cost of the occasional speculative approach).

The participants are usually competent enough to self-assign in a way that skills match the challenge questions well, thus we recommend to focus on conveying the challenge questions as clearly and precisely as possible, while avoiding suggestions of specific approaches.

We recommend around 10 PowerPoint slides in 16:9 format.



Challenge Owner deliverables

The list below gives an overview of Challenge Owner actions and requirements throughout the DSG process.

Pre-event:

- Agree to the terms and conditions of participating in the DSG
- Work with the chosen PI to produce the challenge title, a short description and long description
- Assess the “sensitivity tier” of the data
- Prepare and transfer suitable data to the Turing environment
- Prepare the Monday sales pitch presentation
- Choose a representative to join and get involved in the week

During the event:

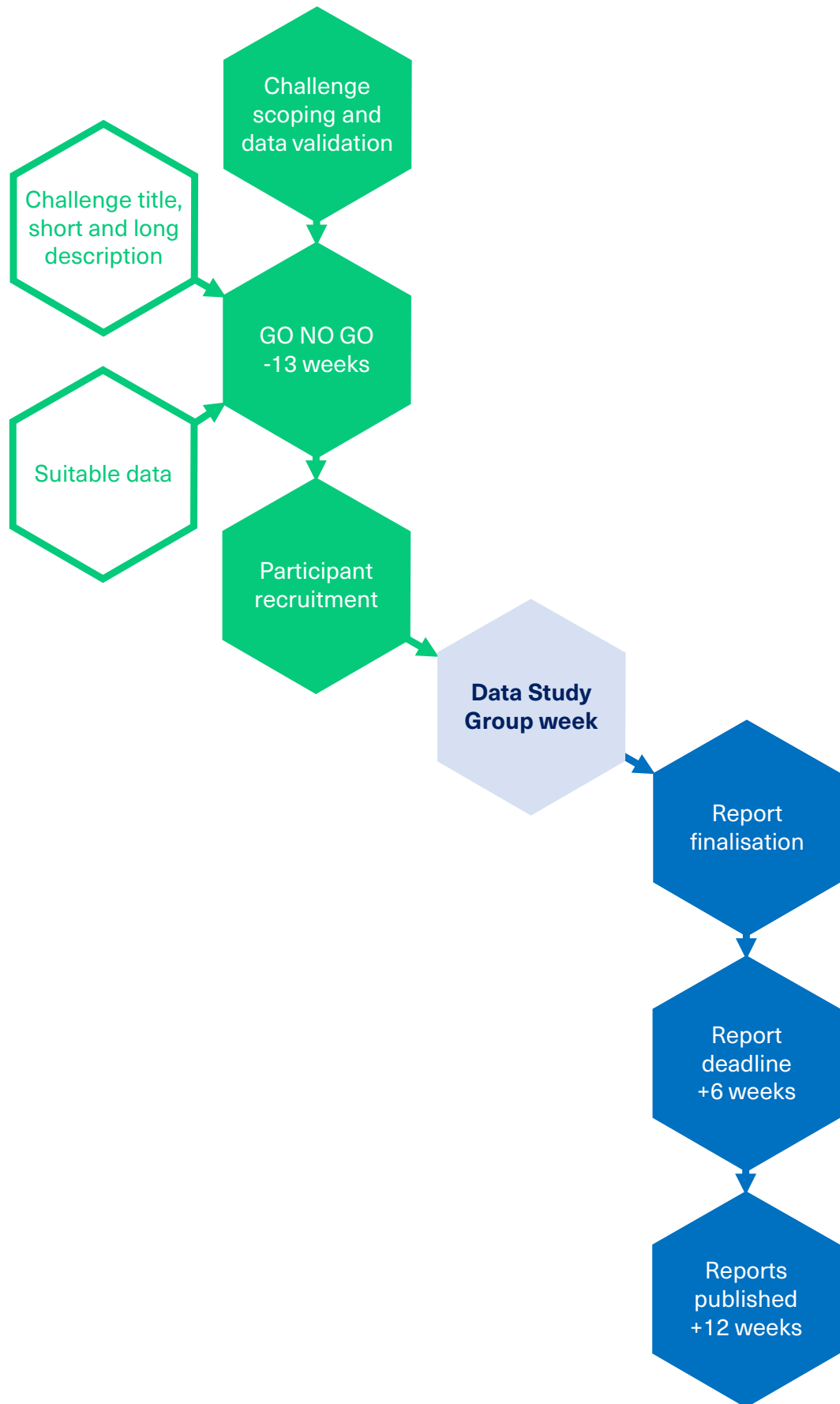
- Attend and deliver the Monday morning presentation
- Answer any questions relating to your challenge
- Attend the Friday presentation
- Ideally try and have a technical member of your team on site to answer any data scientific questions that participants may have, but to also get involved in the group and learn from their teammates

Post-event:

- Engage in the declassification of outputs where necessary
- Review the final report for any commercially sensitive items
- Engage with the Turing for further research (optional)



Challenge cycle timeline



Data security

Each Data Study Group challenge will have its own secure Data Safe Haven* environment used for the storage and analysis of sensitive data. Each environment is separated by an Azure subscription and is its own entity, completely isolated from that of other challenges.

The Safe Haven can include the following security measures depending on data classification†:

- No public internet
- Two factor authentications
- Cut and paste restrictions
- Access to the environment only available from specific sites

For your data to be used in a DSG, the project must be assigned a classification tier which dictates the sensitivity and so the security measures the environment will feature.

As a Challenge Owner you must assign your project a classification tier first, then the PI will assign a tier separately before it is given to a referee for a final assessment. The CO, PI and referee must arrive at the same conclusion before a tier can be successfully assigned.

DSG challenge data classification process

1. The Challenge Owner assigns their project a classification tier
2. The data is moved into a Tier 3 Safe Haven at the Turing
3. The challenge PI then reviews and assigns a tier
4. A referee separately reviews and assigns a tier
5. If all three parties agree to the classification tier, it can be officially assigned. If not, a discussion is had to reach a consensus
6. Environments are then set up for each challenge and the necessary accommodations made surrounding permissions and restrictions depending on the tier ahead of DSG week.

After the Data Study Group, prior to publishing any reports or code from a Tier 1 or above project, a final declassification review is conducted to ensure sensitive information is not published. Here the process is reversed, where the CO, PI and referee review the report (and other outputs), identify items that are deemed sensitive and then discuss how to safely represent said items – delete, anonymise, generalise or obfuscate.

The CO will also have the opportunity to review the report prior to publication, for minor redactions relating to confidentiality.

*For more information on the Data Safe Havens please see turing.ac.uk/research/research-projects/data-safe-havens-cloud

†Definitions of the classification tiers and the flow diagram for assessing the sensitivity can be found in Annex 2

Loading data into the environment

There are two methods of transferring your data to the Safe Haven (in order of preference):

- Microsoft Azure Storage Explorer
- Physically

Microsoft Azure

The Safe Haven is built upon the Microsoft Azure platform. The safest and most convenient way of transferring data from your organisation is to use the Azure Storage Explorer. You do not require log in credentials to access the Azure Storage Explorer.

Before selecting this method, you will need to confirm that you are able to receive a secure email (we use the service egress.com for receiving secure emails).

You will need to send us the public IP address (or range of IP addresses) of the machine that will upload the data. We will ensure that only computers from these IP addresses will be able to connect to our servers.

When we receive the IP address, we will send the person responsible for the data transfer a link via secure email.

Physically

Alternatively, you can bring your data to us on a physical device.

Please let us know when you will come to the Turing to deliver the data.



The role of the Principal Investigator (PI)

A PI will be recruited for each challenge, usually an early career researcher or postdoc from either the Turing or one of our partner universities. The PI will be an experienced data scientist with special expertise relating to your challenge.

The PI will work with the CO on matters that relate to the science and research questions that will be investigated during the DSG. They will help refine your proposal into a challenge suitable for a DSG.

They will also be integral in overseeing the suitability and shape of the data you will provide (e.g. data sensitivity tiers and how that will affect the participant experience during the DSG). The role of the PI will last on average six months per DSG challenge.

They provide instrumental support during all stages of the DSG, including:

- Collaborating with the CO to draw up initial challenge documents such as challenge title and the short and long descriptions.
- Helping to evaluate the appropriateness of data provided to ensure it is challenge-ready before the public challenge announcement.
- Submitting an Ethics Advisory assessment before the DSG to ensure any potential ethical issues have been mitigated.
- Providing report feedback during the DSG week.
- Providing guidance and feedback in finalising the report.
- Leading any follow-on research projects.

Your first meeting with the PI

The primary goal of this first meeting is for both yourself and the PI to sketch out an overview of what the challenge will look like for the DSG.

Your original challenge proposal, along with notes from the challenge assessment, will be shared with the PI prior to this first meeting. However, it would be good to reiterate to the PI what your organisation hopes to achieve from participating in the DSG.

The PI will also have some ideas about how to frame your proposal as a challenge that is both an exciting research question as well as suitable for the timescales of a DSG. It is this balance that should be explored and developed during the meeting.

Do bring a data dictionary to the meeting or send one to the PI beforehand. This is a description of the dataset(s) to be used: the types of information it contains, the dataset's format, structure and size, field headers, what each field means, and how each field relates to other fields.

If permissible, you are strongly encouraged to show the PI a reasonably representative subset of the dataset, so the PI can better understand early on how this data can be used to tackle the challenge.

If you are not a data scientist yourself, it is advisable to also bring someone from your team who is able to discuss at length some of the technical requirements and constraints of your challenge and data. If possible, it is recommended you bring someone who has already worked with the data themselves.

Subsequent meetings will be to iterate on and further refine the challenge description and prepare the data set(s) for the DSG event.

The final report

Reports are a key outcome of the DSG as they provide a tangible output for both organisations and participants, they are publishable by the Turing, and they form the basis for any follow up research work.

During the DSG week, participants will produce reports explaining their work and main findings. On Wednesday or Thursday of the event week, the PI will review the in-progress report and provide feedback to improve it.

Following the DSG event week, there will be a period of report revisions and clear up. Their task will be to collate missing information, complete exposition, finalise the formatting and typesetting, and redact faulty or irreparably incomplete content – all without conducting any further analysis which would be out of scope.

After these iterations, a final editing and review cycle is carried out by DSG editors and the Turing communications team.

If the challenge was Tier 1 or above, a final declassification review is also conducted to ensure that sensitive information is not published.

You can read some of the reports from previous DSGs at: turing.ac.uk/DSG



Attending DSG week

After groups are formed on Monday lunchtime, the real work begins. It is advisable that a representative of your organisation stays to meet the team, give them some more background on the challenge and data, and answer further questions. However, this is not compulsory. Please note that CO travel costs, and accommodation will not be covered by The Alan Turing Institute.

On Friday, the groups will present what they explored during the week. So, we highly recommend the COs attend! You can bring up to four additional guests with you to attend these presentations. You will see all the presentations from all the different groups in this session. If you wish to share your particular presentation more widely, we will be recording it, and we have the facilities to privately livestream.

Get your hands dirty

Depending on your coding skills, joining the group and mucking in with the challenge is an invaluable learning experience. If you are not the coding sort, then we highly recommend that if you have a data scientist in your team you send them to join in.

Not only will they be exposed to new techniques and ideas from the participants, they will also act as the domain knowledge expert for the challenge and can answer all the questions about the data from participants. When the event is over, they will return to your organisation with this knowledge.

Due to the intense nature of the event, we do ask you to be mindful of the following:

- A lot of preparation is done prior to the week to ensure that teams form quickly and start working through their ideas as soon as possible.
- When your colleagues or guests come, by all means feel free to ask the participants questions, however it is vitally important that they do not direct participant avenues of investigation, as this can reset thought processes and group dynamics, which can lead to sub optimal output. Please remember this is not consultancy.

Additional guests

Throughout the week you may wish to invite other members of your organisation to see the progression being made or support the group with additional experts. We welcome your colleagues but do request you inform the Data Study Group team at least two weeks before the event about who will attend and when.

Guests who are not part of your organisation will be required to return a signed non-disclosure agreement (NDA) from their organisation prior to the event. To minimise disruption to the groups we will plan these visits into the individual groups' schedules.

From our experience, a visit on Wednesday morning before lunch is the optimal time to visit – far enough into the challenge to give you a meaningful update, but still enough time to factor in any comments your team may have.

However, it is imperative that you do not tell the participants what to do. They are not employees and have volunteered their time to participate in the DSG. The challenge description is their guide, but participants are free to investigate what they choose, within the bounds of this framework. This is where the serendipity happens.

Code of Conduct

A condition for our participants joining is that they adhere to the DSG Code of Conduct and this extends to all representatives from your organisation who attend the week.

Please find the full version of our Code of Conduct here: bit.ly/2YrEPNx

- Be respectful to others. Do not engage in homophobic, racist, transphobic, ageist, ableist, sexist, or otherwise exclusionary behaviour.
- Use welcoming and inclusive language. Exclusionary comments or jokes, threats or violent language are not acceptable. Do not address others in an angry, intimidating, or demeaning manner.

Be considerate of the ways the words you choose may impact others. Be patient and respectful of the fact that English is a second (or third or fourth!) language for some participants.
- Do not harass people. Harassment includes unwanted physical contact, sexual attention, or repeated social contact. Know that consent is explicit, conscious and continuous—not implied. If you are unsure whether your behaviour towards another person is welcome, ask them. If someone tells you to stop, do so.
- Respect the privacy and safety of others. Do not take photographs of others without their permission. Note that posting (or threatening to post) personally identifying information of others without their consent (“doxing”) is a form of harassment.
- Be considerate of others’ participation. Everyone should have an opportunity to be heard. In update sessions, please keep comments succinct so as to allow maximum engagement by all participants. Do not interrupt others on the basis of disagreement; hold such comments until they have finished speaking.
- Don’t be a bystander. If you see something inappropriate happening, speak up. If you don’t feel comfortable intervening but feel someone should, please feel free to ask a member of the Code of Conduct response team for support.
- As an overriding general rule, please be intentional in your actions and humble in your mistakes.

Follow-on collaboration

It is important to be aware that as well as training the next generation of data scientists, the DSG serves to kick start research collaborations between the Turing, its partner universities, and industry, government and the third sector.

It is our aim that we try to continue as many of the projects that we have started at the DSG in some form of follow-on work. This could be from writing a paper to extending the research, working closely with the CO continuing to investigate the results from the DSG more closely, to forming longer term partnerships.

While working with the PI in the preparation of the challenges, this will be one of the things that they will consider – how this work will lead on to something bigger.

Data Study Groups have been a great launchpad for follow-on research in the past, please follow the link to browse a selection of research that has in some way come out of a Data Study Group.

turing.ac.uk/DSG

Data Study Group contacts



Sebastian Vollmer – Data Study Groups Director

Sebastian invented the Data Study Groups with the aim of fostering challenge-driven research and translational activities as an inclusive community effort.

Sebastian is a Turing Fellow, an associate Director on the Health programme at the Turing, and an Associate Professor for Mathematics and Statistics at the University of Warwick. His work has two pillars: foundations of Bayesian inference and Monte Carlo methods and applications of data science.

He enjoys his free time with his family and likes board games and discussing politics.



Jules Manser - Data Study Group Project Manager

Jules is a member of the Partnerships team at the Turing and has played a lead role in delivering Data Study Groups, working closely with Sebastian since their inception in 2016.

Alternatively, please contact your Turing partnership development contact to see if the Data Study Group is a viable route for your challenge.

FAQs

When is the DSG?

DSG is held around three times a year at the Turing offices in London. A single DSG event usually hosts around 5 – 6 challenges.

How is the DSG week structured?

Day 1: Challenges are presented by Challenge Owners in the morning, participants self-select which challenge they want to participate in after lunch, then begin to brainstorm.

Day 2-4: Brainstorming, modelling and problem solving.

Day 5: Progress and recommended routes forward are presented.

Is the DSG a cheap consultant?

No. Whilst the CO dictates the scope of the challenge, it is up to the PI to shape the challenge into an interesting research question. The resultant question should provide a framework for the participants to explore during the week.

Participants are offered the freedom to investigate the data as they wish and use the challenge description as a guide. What is more, the report findings and code will be published on the website, available for all to share and learn from.

What about Intellectual Property rights?

Any intellectual property arising from the DSG will be owned by the Turing. All of the CO's background IP and that of our researchers, remains with the inventor. For 3.5 days of work we do not expect any patents to be registered. Remember, all results will be published, and any code developed will be made available under permissive open source license.

What is a facilitator?

A facilitator is chosen from the pool of applicants. Their job is to manage the group, specifically people dynamics. They will not be heavily involved in the data analysis but will have had a preview of the dataset before the event week to help support the other participants with familiarising themselves with the data.

How are participants selected?

Applicants apply directly online via the Turing website. The opportunity is circulated around the Turing network, to our university partners and further afield through external communications encouraging scientists/researchers to apply.

All applications are reviewed by at least two members of an experienced selection team and scored on factors such as technical and collaborative ability.

Annex 1

Example of long challenge description

Accenture: 'Fairness in algorithmic decision-making'

Background

Accenture provides end-to-end services in strategy, consulting, digital, technology, and operations. With expertise across more than 40 industries and a global presence, we are helping business transform into the digital world.

Accenture Applied Intelligence is a practice within Accenture's Digital arm working with clients on intelligent technologies and driving analytical advancements in industry. As part of this, we and our clients see the importance of using Artificial Intelligence in a responsible way. We want to help clients use AI to the fullest in a confident and responsible manner.

As AI becomes an integral part of business processes, we begin to see both the positive and negative aspects of this technology. On the positive, processes that were previously subjective, based on human inputs (e.g. loan approvals, claims disputes) can be processed in a systematic and data-rich manner.

On the negative, we realise that the data used to train these algorithms, and sometimes the algorithms themselves, can result in discriminatory applications of the output. As a result, many are calling for algorithms to ensure fairness.

Transparency in algorithmic decision-making is critical to the successful implementation of AI, as it enables trust between consumers and organisations. Being clear about what the definition of 'fairness' is, and having a systemic way of assessing multiple viable definitions of fairness to select for the best outcome, is one method of improving transparency. In the banking industry, algorithms are being used to create more data-driven outcomes, but the critical decisions being made require that the methods used are designed with care.

The challenge

In this challenge we wish to examine the concept of algorithmic fairness applied in a business setting, with a focus on the financial services industry.

In the financial services industry, algorithms significantly impact the services customers are able to receive, such as credit cards or mortgages. We want to help financial institutions better understand the concept of fairness (as required e.g. by the GDPR) by use of a framework and/or a set of metrics. In the case that these concepts remain nebulous it will be challenging for organisations to implement them or consumers to benefit.

We understand and emphasise that fairness cannot be exhaustively codified, but the goals are pragmatic: how can we help industry data scientists think carefully about these issues? We aim to use multiple definitions of fairness, codify and quantify them for the financial services industry, and create an understandable rubric.

The goal for this project is for individuals utilising an algorithm for potentially sensitive outcomes to be able to illustrate the trade-offs between different definitions of fairness and select the best definition to suit the algorithmic implications. Such work may then be applied to create an Algorithmic Impact Assessment, providing clarity into decisions.

The data

We will be utilising a publicly available dataset on credit risk decision-making from the UCI Repository. This dataset classifies people described by a set of attributes as good or bad credit risks.

Included in the data are sensitive variables, including gender, employment status - including type of employment, marital status, and residency status (foreign or not).

Goals of the Data Study Group week

- The key aim would be towards creating a reusable rubric for measuring fairness of algorithmic outputs, highlighting the strengths and weaknesses of each.
- A starting point for this work is Dr Arvind Narayanan's tutorial at the FATML conference, 21 fairness definitions. There are a large number of different metrics and to the best of our knowledge, no comprehensive work has considered these jointly. We are looking to implement as many of these as possible in order to gain insight and intuition into what they capture, how they interplay and what the relative merits are for each.
- We are further interested in any synthesis of these existing metrics which may both capture ideas of fairness more fairly and be more practical for end users.
- We also want to consider a broader question: what notable gaps are left from use of these metrics? In which contexts are they more appropriate, and how context sensitive are the remaining problems?
- Understanding how existing theory might be tailored towards different industry sectors. We are particularly interested in the financial sector, where we will apply the results to a banking dataset scoring credit risk. We aim to measure the fairness of scoring each relevant individual as 'high risk'.
- Finally, we look to create a general, easy to apply list of definitions and metrics so data scientists and industry SMEs can choose the measurement that is the most relevant to their data, model, and industrial context.

What we're looking for

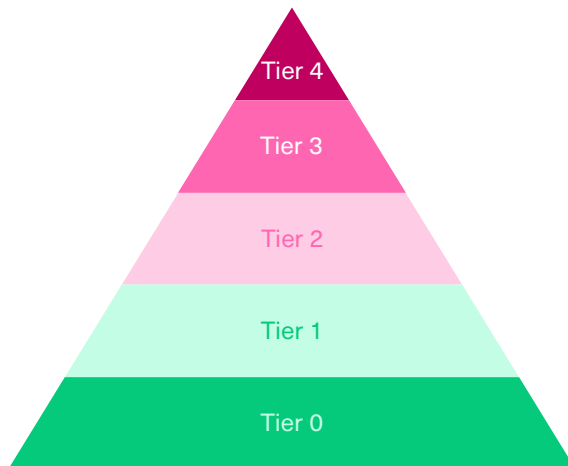
We are seeking quantitative scientists as well as ethicists and philosophers familiar with different definitions of societal and individual fairness. We will ask the group to consider the social context of the data, be creative in the development of metrics, and work with Accenture's subject matter experts in Responsible AI and financial services to capture industry nuances.

If we have additional time, we may explore the outcomes of multiple algorithms to illustrate how fairness measurements may change for different algorithmic output across the different fairness criteria.

Annex 2

Sensitivity tiers and assessment

Turing Data Safe Haven tiers and flow



Our model goes from Tier 0 – publicly available, open information – to Tier 4 – personal data where disclosure poses a substantial risk to safety.

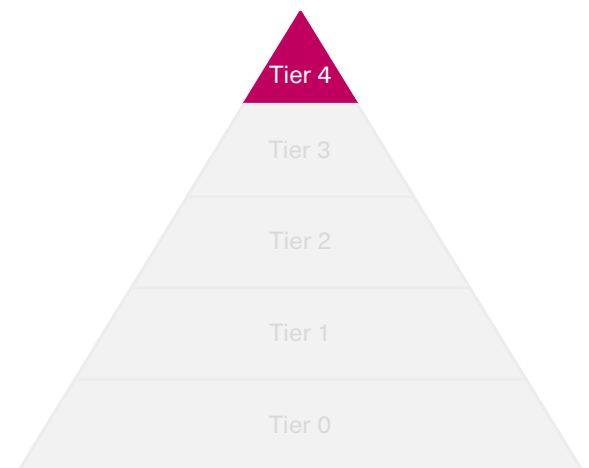
This guidance document should give you an idea of how to classify your data for the Data Study Group.

It should be referenced in conjunction with the classification flowchart.

Tier 4 environments are for:

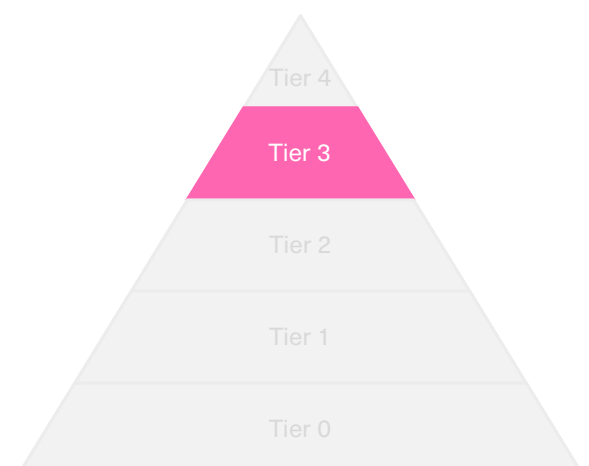
- Personal data where disclosure poses a substantial threat to safety, security or health
- Commercial or governmental data which could be subject to attack by sophisticated, well-resourced and determined actors such as nation states

Tier 4 data is **not** appropriate for Data Study Group use.



Tier 3 environments are for:

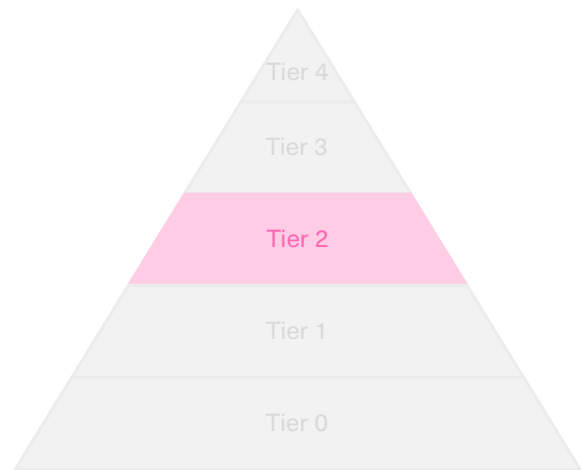
- Pseudonymised or synthetic data where confidence in the quality of anonymisation is weak
- Commercial data which is sensitive
- Commercial or governmental data which could be subject to attack by attackers with bounded capabilities such as hackers



Tier 2 environments are for:

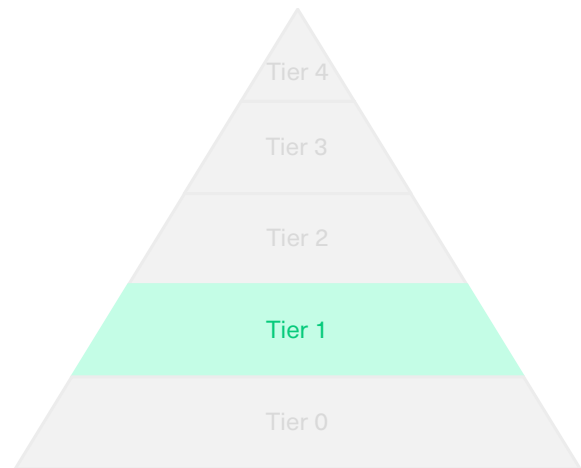
- Pseudonymised or synthetic data where confidence in the quality of anonymisation is strong
- Commercial data where risks from disclosure are low

Tier 2 data should be very unlikely to be subject to targeted attack.



Tier 1 environments are for:

- Data intended for eventual but not immediate publication
- Datasets where the only risks of disclosure are to the researchers' competitive advantage
- Pseudonymised or synthetic data where confidence in the quality of anonymisation is absolute
- Commercial information where the consequences of disclosure are so low as to be trivial

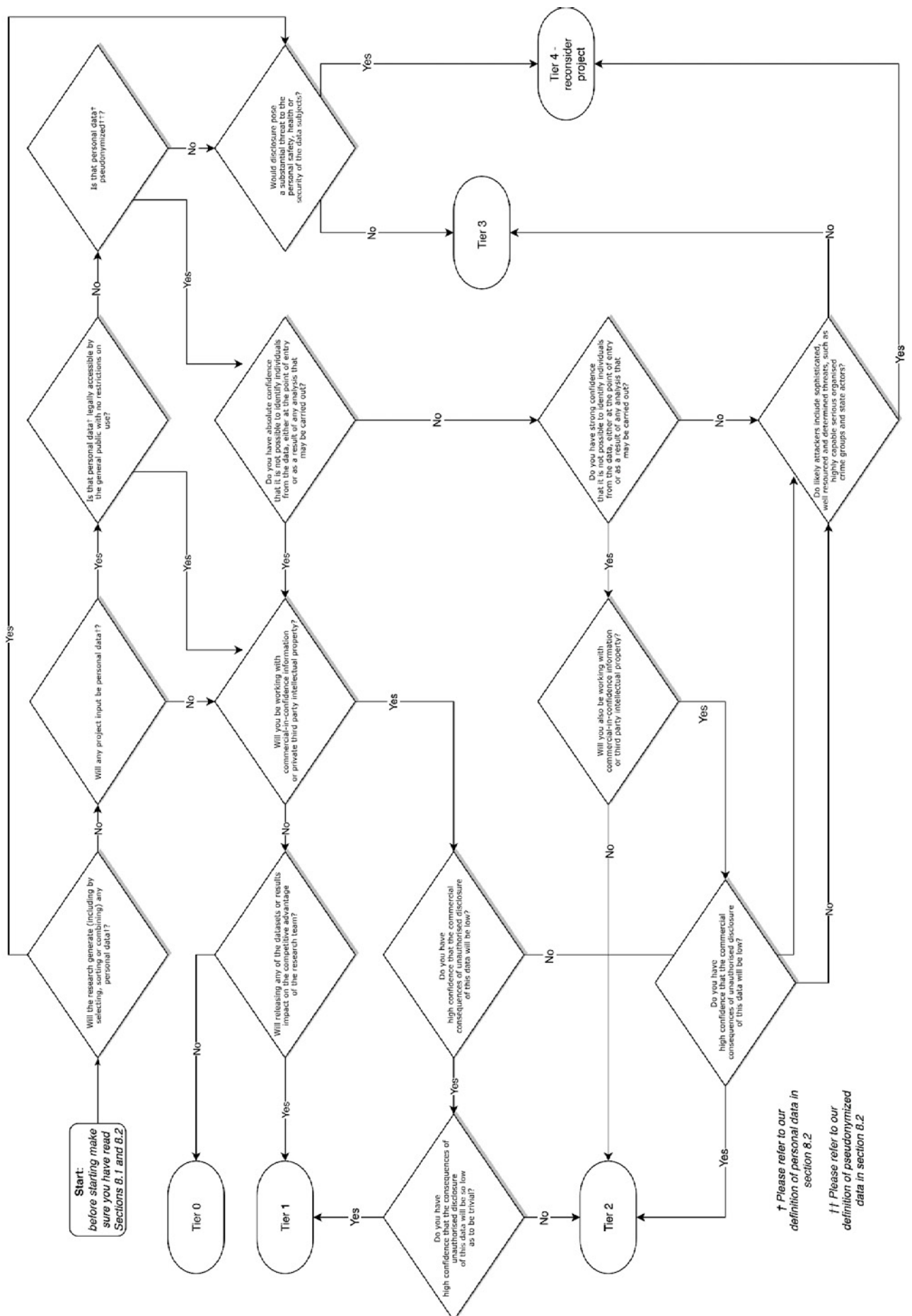


Tier 0 environments are for:

- Publicly available, openly published information
- Data which is intended for immediate publication



How to assess the tier of the projects





turing.ac.uk
@turinginst