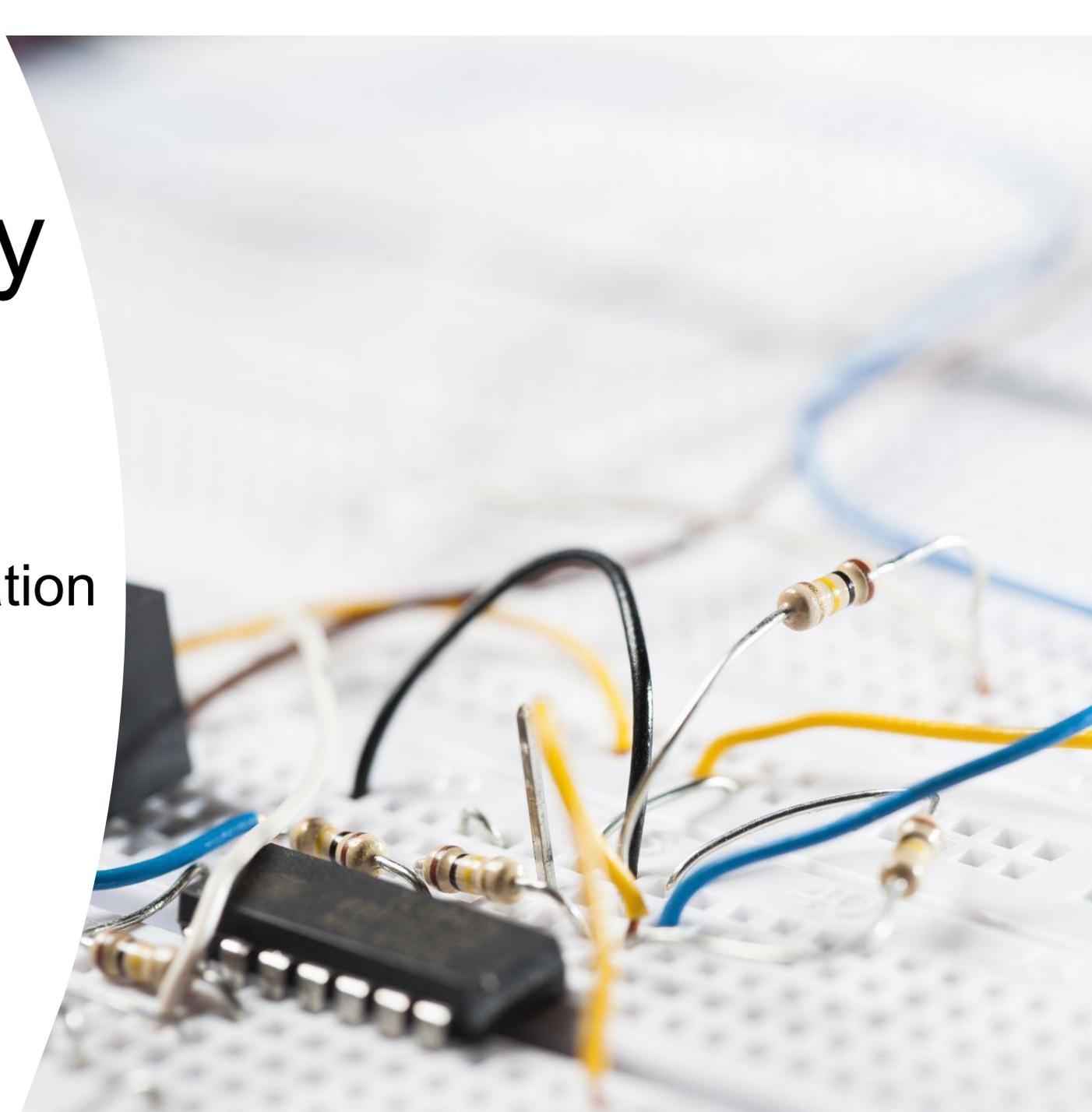# TRM Compulsory Module

Misleading Statistics & Visualisation
Writing Introduction
Writing Methodology

Dr. Mingshan Jia
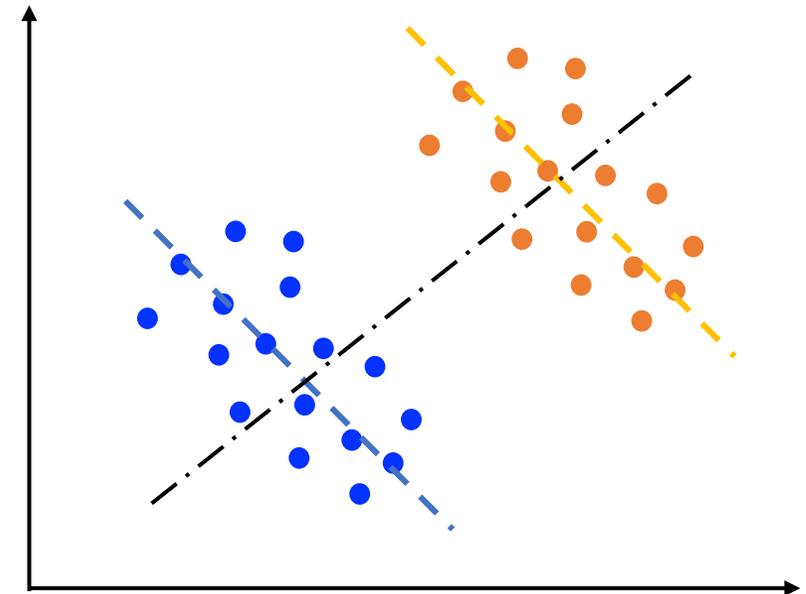
# 1. Misleading Statistics & Visualisation

# Agenda

- Misleading Statistics
  - Simpson's paradox
  - Berkson's bias
  - Inspection paradox
- Misleading visualisations
- Activity

# Simpson's Paradox

Simpson's paradox is a statistical paradox where it's possible to draw two opposite conclusions from the same data depending on how we group data.

- Aggregated data can sometimes hide the truth.

- The direction of a trend can reverse when data is divided into subgroups, and vice versa.

- This usually happens due to confounding variables.

# Simpson's Paradox: Admission Bias

A popular Simpson's paradox is the admission to graduate study at the University of California, Berkeley for the fall 1973 quarter.

- Overall data for university admissions: Bias against women

| | Men | | | Women | | |
|---|---|---|---|---|---|---|
| | Applicants | Admitted | Admission ratio | Applicants | Admitted | Admission ratio |
| | 100 | 40 | 40% | 100 | 30 | 30% |

- By faculty: Bias against men

| Faculty | Men | | | Women | | |
|---|---|---|---|---|---|---|
| | Applicants | Admitted | Admission ratio | Applicants | Admitted | Admission ratio |
| Engineering | 80 | 38 | 48% | 20 | 14 | 70% |
| Arts & humanity | 20 | 2 | 10% | 80 | 16 | 20% |

Source

# Simpson's Paradox: Who is a better shooter?

- Overall Shooting Accuracy: LBJ and TD are Equally good.

| Player | FGM | FGA | FG% |
|--------|-----|-----|-----|
| LeBron James | 15,842 | 31,285 | 50.60% |
| Tim Duncan | 10,285 | 20,334 | 50.60% |

- By 2 and 3 pointers: Lebron James is a better shooter.

| Player | 2 Pointers | | | 3 Pointers | | |
|--------|------|------|------|------|------|------|
| | 2PM | 2PA | 2P% | 3PM | 3PA | 3P% |
| LeBron James | 13,200 | 23,762 | 55.60% | 2,622 | 7,523 | 34.90% |
| Tim Duncan | 10,255 | 20,166 | 50.90% | 30 | 168 | 17.90% |

# Berkson's Bias
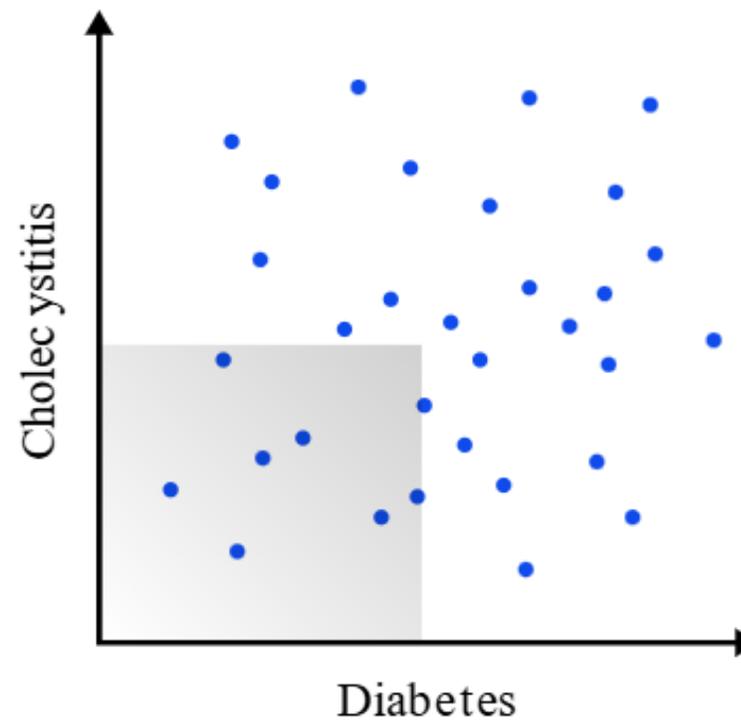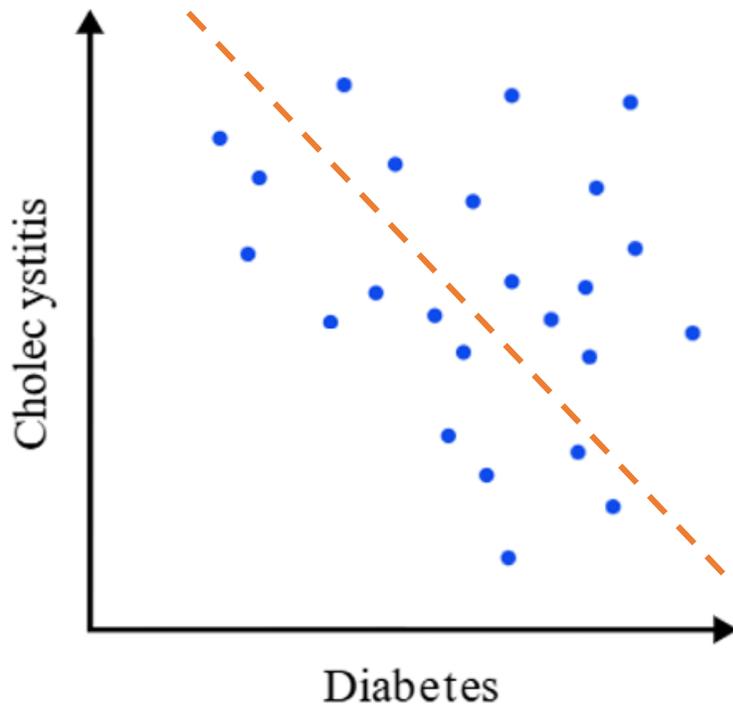
Berkson's bias or Selection Bias is the counterintuitive idea that events that seem to be correlated are actually not correlated or correlated in the inverse direction.

- It arises from non-random sampling, often due to how participants are selected.

- It tends to produce spurious correlations (often negative) between independent variables.

- It occurs when the probability of inclusion in the sample is related to the variables being studied.
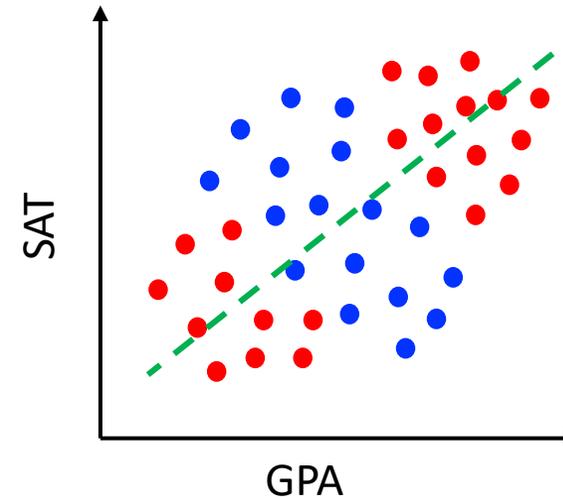


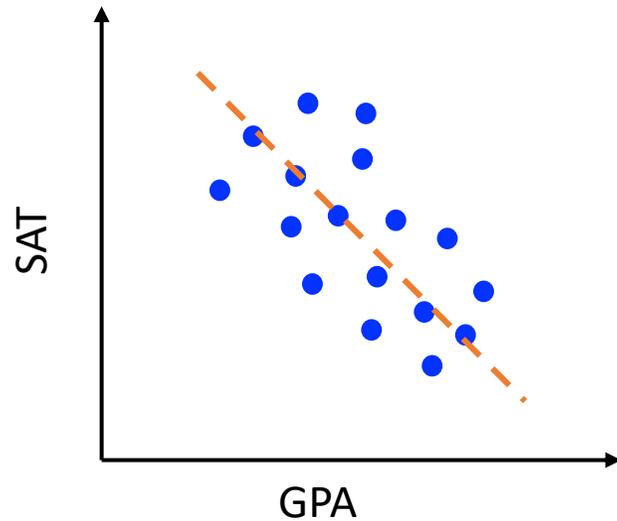The selected sample

The entire population

# Berkson's Bias: Disease Correlation?

For example, two studied diseases are diabetes and cholecystitis, and based on the collected data from in-hospital patients, we observe a negative correlation between the two.

# Berkson's Bias: College Admission

At a university, we observe a negative correlation between GPA and SAT score among newly enrolled students.

# Inspection Paradox

The Inspection Paradox describes situations where the experience of individuals is different from the overall system average — due to the way observations are sampled.

- The paradox is caused from biased sampling:
  - Size-biased sampling
  - Length-biased sampling

- The difference is essentially different ways of calculating the average
  - a simple average
  - a weighted average

# Inspection Paradox: Flight Occupancy

- Passengers always find their flights full, but the airlines always complain about low occupancy rates.



Inspection Paradox: Perceived vs Actual Flight Occupancy

# Inspection Paradox: Flight Occupancy

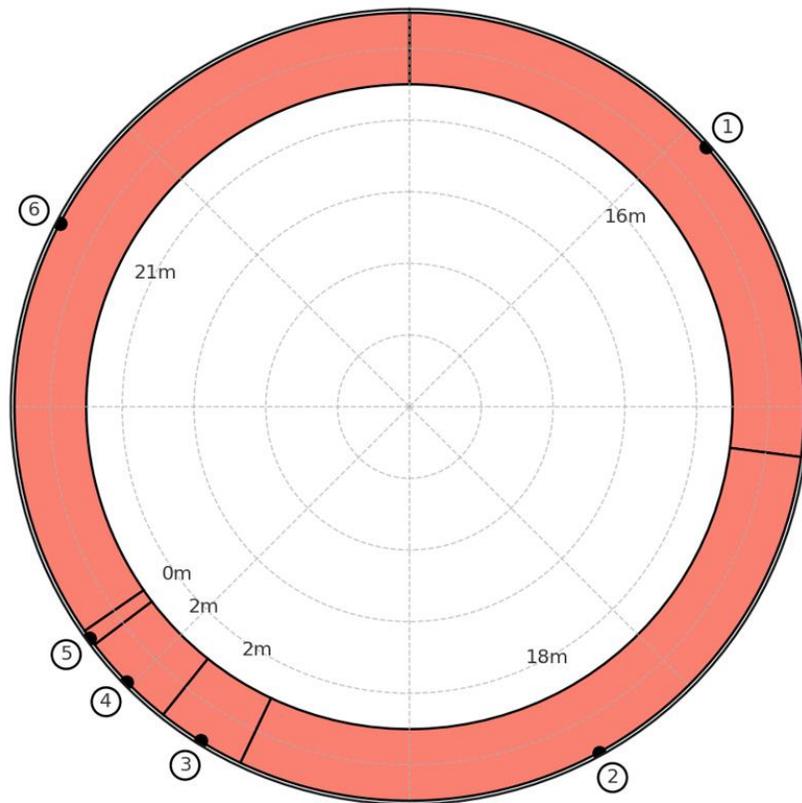- Passengers always find their flights full, but the airlines always complain about low occupancy rates.

| Flight No. | Number of Passengers |
|:---:|:---:|
| A | 10 |
| B | 20 |
| C | 90 |

- System average (what the airline sees):
  - (10 + 20 + 90) / 3 = 40
- Passenger-Experienced Average (what a random passenger sees):
  - (10/120) * 10 + (20/120) * 20 + (90/120) * 90 =  71.7

# Inspection Paradox: Bus Waiting Time

- **Bus Waiting Paradox**: The experienced bus waiting time is longer than the average bus interval.



6 Bus Arrivals in One Hour

- According to bus scheduling system:
  - Average interval: 60/6 = 10 mins
  - System average waiting time: 10/2 = **5 mins**

- According to experienced waiting time:
  - Passenger-experienced average waiting time: **8 mins**

$$\sum_{i=1}^{n} \left( \frac{l_i}{\sum_{j=1}^{n} l_j} \cdot \frac{l_i}{2} \right)$$

# Misleading Visualisation

**Truncated Axes**

Truncating axes can make minor differences look dramatic.

**3D Distortion**

3D visualisation can distort relative size by exaggerating objects in the foreground.

**Inconsistent Scale**

Multiple groups are placed side by side but not scaled equally.

# Misleading Visualisation

- **Truncated Y-Axis**



Fox Business broadcasted on August 1, 2012

- **Distortion of 3D Visualisation**



WWDC 2008, June 9, 2008

Research Paper: Truncating Bar Graphs Persistently Misleads Viewers

- Misleading Visualisation
  --- Inconsistent Scale



The number of artists generating at least $1M, $100K and $10K, has **nearly tripled** since 2017

Source

# Activity

1. Group Discussion
   Discuss with your group members:
   - Have you ever encountered misleading statistics or visualisations?
   - If yes, what made them misleading?

2. Reflect on Your Own Experiment
   - What results do you expect to obtain?
   - What type of visualization do you plan to present your results?
   - Is there any potential bias or misleading presentation you should avoid?

# 2. Writing Introduction

# Agenda

- The Audience of a Research Paper

- Key Goals of the Introduction

- Typical Structure

- Example

# The Audience of a Research Paper

Academia:

- Students and early-career researchers

- Peer researchers and scholars

- Reviewers

Outside of Academia:

- Industry Practitioners

- Funding bodies

# Key Goals of Writing the Introduction

Communicate the gist of the research to an audience

Let readers understand the broad research topic and its significance

[General]

Let readers understand the contributions (theoretical) and findings (empirical)

[Specific]

- Make the readers feel interested
- Make the readers feel they've learned something

# Typical Structure of the Introduction

- Introduce the broad research area and show its significance (context).

- Discuss the limitations of existing works (research gap).

- Introduce the proposed approach (theoretical contribution / novelty).

- Present key findings (empirical evidence).

- Summarise key contributions.

# Typical Structure: Example 1

**FineTuning Text-to-Image Diffusion Models for Fairness (ICLR 2024)**

**Hypothesis**: *Biases in T2I diffusion models can be mitigated in a controllable manner through finetuning.*

**Research Question**: How *can we fine-tune text-to-image diffusion models to mitigate bias in a controllable manner, ensuring that generated images align with a user-defined target distribution for attributes such as gender, race, and age ?*

# Typical Structure: Example 1

- Paragraph 1: The popularity of T2I diffusion models.

## 1 INTRODUCTION

Text-to-image (T2I) diffusion models (Nichol et al., 2021; Saharia et al., 2022) have witnessed an accelerated adoption by corporations and individuals alike. The scale of images generated by these models is staggering. To provide a perspective, DALL-E 2 (Ramesh et al., 2022) is used by over one million users (Bastian, 2022), while the open-access Stable Diffusion (SD) (Rombach et al., 2022) is utilized by over ten million users (Fatunde & Tse, 2022). These figures will continue to rise.

- Paragraph 2: Bias issue of T2I diffusion models, and limitations of existing approaches.

However, this influx of content from diffusion models into society underscores an urgent need to address their biases. Recent scholarship has demonstrated the existence of occupational biases (Seshadri et al., 2023), a concentrated spectrum of skin tones (Cho et al., 2023), and stereotypical associations (Schramowski et al., 2023) within diffusion models. While existing diffusion debiasing methods (Friedrich et al., 2023; Bansal et al., 2022; Chuang et al., 2023; Orgad et al., 2023) offer some advantages, such as being lightweight, they struggle to adapt to a wide range of prompts. Furthermore, they only approximately remove the biased associations but do not offer a way to *control* the distribution of generated images. This is concerning because perceptions of fairness can vary across specific issues and contexts; absolute equality might not always be the ideal outcome.

# Typical Structure: Example 1

- Paragraph 3: Proposed approach (Part 1: DAL)

We frame fairness as a distributional alignment problem, where the objective is to align particular attributes of the generated images, such as gender, with a user-defined target distribution. Our solution consists of two main technical contributions. First, we design a loss function that steers the generated images towards the desired distribution while preserving image semantics. A key component is the distributional alignment loss (DAL). For a batch of generated images, DAL uses pre-trained classifiers to estimate class probabilities (*e.g.*, male and female probabilities) and dynamically generates target classes that match the target distribution and have the minimum transport distance. To preserve image semantics, we regularize CLIP (Radford et al., 2021) and DINO (Oquab et al., 2023) similarities between images generated by the original and finetuned models.

- Paragraph 4: Proposed approach (Part 2: Adjusted DFT)

Second, we propose adjusted direct finetuning of diffusion models, adjusted DFT for short, illustrated in Fig. 2. While most diffusion finetuning methods (Gal et al., 2023; Zhang & Agrawala, 2023; Brooks et al., 2023; Dai et al., 2023) use the same denoising diffusion loss from pre-training, DFT aims to directly finetune the diffusion model's sampling process to minimize any loss defined on the generated images, such as ours. However, we show the exact gradient of the sampling process has exploding norm and variance, rendering the naive DFT ineffective (illustrated in Fig. 1). Adjusted DFT leverages an adjusted gradient to overcome these issues. It opens venues for more

# Typical Structure: Example 1

- Paragraph 5: Empirical Findings (Part 1: Effective debiasing; Effective finetuning )

Empirically, we show our method markedly reduces gender, racial, and their intersectional biases for occupational prompts. The debiasing is effective even for prompts with unseen styles and contexts, such as "*A philosopher reading. Oil painting*" and "*bartender at willard intercontinental makes mint julep*" (Fig. 3). Our method is adaptable to any component of the diffusion model being finetuned. Ablation study shows that finetuning the text encoder while keeping the U-Net unchanged hits a sweet spot that effectively mitigates biases and lessens potential negative effects on image quality. Surprisingly, finetuning as few as five soft tokens as a prompt prefix is able to largely reduces gender bias, demonstrating the effectiveness of soft prompt tuning (Lester et al., 2021; Li & Liang, 2021) for fairness. These results underscore the robustness of our method and the efficacy of debiasing T2I diffusion models by finetuning their language understanding components.

- Paragraph 6: Empirical Findings (Part 2: User-defined distribution; Debiasing multiple concepts)

A salient feature of our method is its flexibility, allowing users to specify the desired target distribution. For example, we can effectively adjust the age distribution to achieve a 75% young and 25% old ratio (Fig. 4) while *simultaneously debiasing gender and race* (Tab. 5). We also show the scalability of our method. It can debias multiple concepts at once, such as occupations, sports, and personal descriptors, by expanding the set of prompts used for finetuning.

- Paragraph 7: Conclusion

# Typical Structure: Example 2

Measuring Quadrangle Formation in ComplexNetworks (TNSE 2021)

**Hypothesis**: *Novel coefficients that quantify the formation of quadrangles are more suitable to describe and study quadrangle-rich networks (e.g., biological, ecological, and infrastructure networks).*

**Research Question**: *How can we effectively measure the formation of quadrangles in complex networks, and what insights can these new measurements bring to network analysis?*

# Typical Structure: Example 2

- Paragraph 1: The importance of studying local structure.

- Paragraph 2: The classic metric for studying triangle formation.

- Paragraph 3: A recent novel perspective on studying triangle formation

- Paragraph 4: The importance of quadrangles in real-world networks

- Paragraph 5: Proposed approach (Part 1)

- Paragraph 6: Proposed approach (Part 2: Extend to bipartite networks)

- Paragraph 7: Proposed approach (Part 2: Extend to weighted networks)

- Paragraph 8: Empirical Findings (Part 1: Effective in network description)

- Paragraph 9: Empirical Findings (Part 2:  Effective in network analysis)

- Paragraph 10: Summary

```
%\textcolor{RubineRed}{Para1: the importance of
studying local structure}\\
\IEEEPARstart{C}{omplex} systems across various
domains, such as biology, ecology, physics and
social science, can be modelled as networks that
abstract the interactions between system's
components \cite{barabasi2016network,
newman2018networks, musial2013social}. Different
from a simple grid graph or a line graph for image
or text modelling respectively, the complexity of
networks comes from their intricate topological
structures. Therefore, the study of network
structure, especially local structure, underlies a
number of representative and analytical
applications such as representation learning of
graphs
\cite{hamilton2017representation,grover2016node2vec
}, node-type classification \cite{bhagat2011node,
kipf2016semi}, link prediction \cite{gao2015link,
kovacs2019network} and anomaly detection
\cite{noble2003graph, akoglu2015graph}.

%\textcolor{RubineRed}{Para2: clustering-co}
One fundamental and classic statistical metric to
assess the local structure of complex networks is
the \textit{local clustering coefficient}
\cite{watts1998collective, fagiolo2007clustering}.
It is defined as the percentage of the number of
triangles formed with a focal node to the number
of triangles that the focal node could form with
all its neighbours. Note that the focal node here
serves as the centre node in an open triad (the
middle of a length-2 path). Since many of the real-
world networks are triangle-rich, the clustering
coefficient --- a measure of triangle formation --
- has become a standard metric to describe
networks. It has also been used in numerous
applications such as malware detection
\cite{lee2018automatic}, language learning
\cite{goldstein2014influence} and structural role
discovery \cite{henderson2012rolx}.

%\textcolor{RubineRed}{Para3: closure-co}
A recent study has proposed another interesting
measure of triangle formation, i.e., the
\textit{local closure coefficient}
```

# 1 INTRODUCTION

COMPLEX systems across various domains, such as biology, ecology, physics and social science, can be modelled as networks that abstract the interactions between system's components [1], [2], [3]. Different from a simple grid graph or a line graph for image or text modelling respectively, the complexity of networks comes from their intricate topological structures. Therefore, the study of network structure, especially local structure, underlies a number of representative and analytical applications such as representation learning of graphs [4], [5], node-type classification [6], [7], link prediction [8], [9] and anomaly detection [10], [11].

One fundamental and classic statistical metric to assess the local structure of complex networks is the *local clustering coefficient* [12], [13]. It is defined as the percentage of the number of triangles formed with a focal node to the number of triangles that the focal node could form with all its neighbours. Note that the focal node here serves as the centre node in an open triad (the middle of a length-2 path). Since many of the real-world networks are triangle-rich, the clustering coefficient — a measure of triangle formation — has become a standard metric to describe networks. It has also been used in numerous applications such as malware detection [14], language learning [15] and structural role discovery [16].

A recent study has proposed another interesting measure of triangle formation, i.e., the *local closure coefficient* [17]. With the focal node as the end node of an open triad (the head of a length-2 path), it is quantified as the percentage of twice the number of triangles containing the focal node to the

• M. Jia, B. Gabrys and K. Musial are with the School of Computer Science, University of Technology Sydney, Ultimo NSW 2007, Australia. E-mail: mingshan.jia@student.uts.edu.au, {bogdan.gabrys, katarzyna.musial-gabrys}@uts.edu.au
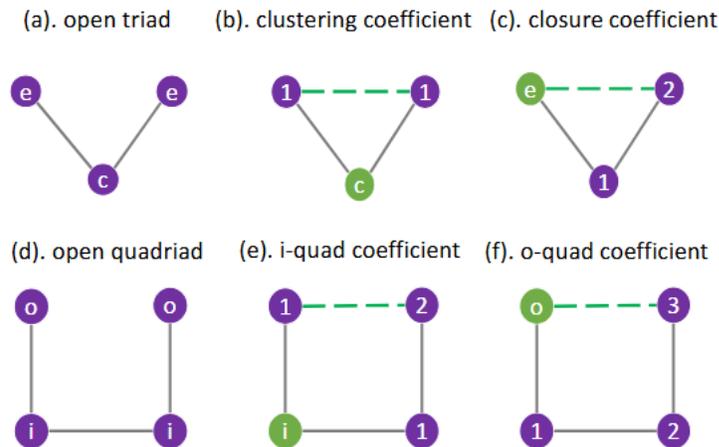
Fig. 1: The i-quad coefficient and the o-quad coefficient in comparison with the clustering coefficient and the closure coefficient. Letters $c$, $e$, $i$ and $o$ denote centre node, end node, inner node and outer node respectively. Node in green colour is the focal node in each subfigure. Number on node indicates the node's distance from the focal node in the open triad or the open quadriad, which might be closed by an edge in dotted green line style.

number of all length-2 paths starting from the focal node. Specifically, the classic local clustering coefficient measures the extent to which the 1-hop neighbours of a given node connect to each other, while the local closure coefficient measures the extent to which the 2-hop neighbours of a given node connect to the given node itself. This new metric has been proven to be a useful feature in network analysis tasks such as community detection and link prediction [17].

In many types of networks, however, quadrangles appear at a much higher frequency than triangles, and thus become the most dominant motifs [18]. For instance, in gene

# Tips for Writing the Introduction

- Use easy-to-understand language

- Ensure smooth and logical transitions between paragraphs

- Illustrate with examples

- Include a figure
  - Illustrate the motivation
  - Explain the core idea
  - Visually present the results (especially when the results are intuitive or visually appealing)

# Activity

1. Choose a research paper of your interest, and analyse how its introduction section is written.

- What is the broad research topic, and why it is worth researching?
- What is the identified gap and how it leads to the proposed approach?
- Is there a figure included? If so, what role does it play?

2. Prepare for your own introduction section:

- Outline the key ideas of each paragraph you plan to write.
- Think of a figure that can help enhance your introduction.

# 3. Writing Methodology

# Agenda

- Where Methodology Sits in a Research Paper
- Key Goals of the Methodology Section
- Typical Structure
- Tips for writing the Methodology
- Homework

# Where Methodology Sits in a Research Paper

Introduction → **Methodology** → Experiment

- Introduce the broad research area and show its significance (context).

- Discuss the limitations of existing works (research gap).

- Introduce the proposed approach (theoretical contribution / novelty).

- Present findings (empirical evidence).

- Summarise key contributions.

# Key Goals of the Methodology Section

The goal of the Methodology section is to present the proposed approach in detail and convince readers the approach should work by design.

- Present the core procedures of the proposed approach.
  - Use pipeline diagrams, flowcharts, algorithms, etc.
  - Allow others to replicate your work

- Justify every design choice
  - Give motivation and justification through discussion

- Demonstrate rigour and credibility
  - Show theoretical foundation and reasoning process.
  - Use definitions, mathematical notations, equations, algorithms, and proofs.

# Structure of the Methodology Section

There is no one-size-fits-all template for the methodology section.

- Starts with a heading called "Method", "Methodology", "Proposed Approach", or the actual method name.

- Then followed by sub-headings that mirror the pipeline or building blocks of the method.

Example 1:
4 Method
- 4.1 Loss Design
- 4.2 Adjusted Direct Finetuning of Diffusion Model's Sampling Process

Example 2:
3 Two Quadrangle Coefficients
- 3.1 I-quad coefficient
- 3.2 O-quad coefficient
- 3.3 Quadrangle coefficients in weighted networks
- 3.4 Computational cost

# Tips for Writing the Methodology Section

## Keep it clear and understandable

- Use easy-to-understand language
- Use logical structures
- Explain equations
- Discuss & Exemplify complex concepts/processes
- Visualise the workflow
- (Present the algorithm)

## Show rigour and credibility

- Provide critical details
- Release code and data
- (Give proofs/derivations)

# Tips: Example

- **Use logical structure:**

### 4.1 LOSS DESIGN

**General case**    For a clearer introduction, we first present the loss design for a general case, which consists of the distributional alignment loss $\mathcal{L}_{\text{align}}$ and the image semantics preserving loss $\mathcal{L}_{\text{img}}$. We

**Adaptation for face-centric attributes**    In this work, we focus on face-centric attributes such as gender, race, and age. We find the following adaptation from the general case yields the best results.

- **Explain equations:**

$$\sigma^* = \operatorname*{argmin}_{\sigma \in \mathcal{S}_N} \sum_{i=1}^{N} |p^{(i)} - u^{(\sigma_i)}|_2, \tag{4}$$

where $\mathcal{S}_N$ denotes all permutations of $[N]$, $\sigma = [\sigma_1, \cdots, \sigma_N]$, and $\sigma_i \in [N]$. Intuitively, $\sigma^*$ finds, in the class probability space, the most efficient modification of the current images to match the target distribution. We construct $\{u^{(i)}\}_{i \in [N]}$ to be i.i.d. samples from the target distribution and

Having detected the cause of the issue, we propose adjusted DFT, which uses an adjusted gradient that sets $A_t = 1$ and $B_t = \mathbf{I}$: $\left(\frac{d\mathbf{z}_0}{d\theta}\right)_{\text{adjusted}} = -\sum_{t=1}^{T} \frac{\partial \boldsymbol{\epsilon}^{(t)}}{\theta}$. It is motivated from the unrolled expression of the reverse process:

$$\mathbf{z}_0 = -\sum_{t=1}^{T} A_t \boldsymbol{\epsilon}_\theta(g_\phi(\mathbb{P}), \mathbf{z}_t, t) + \frac{1}{\sqrt{\bar{\alpha}_T}} \mathbf{z}_T + \sum_{t=2}^{T} \frac{1}{\sqrt{\bar{\alpha}_{t-1}}} \boldsymbol{w}_t, \boldsymbol{w}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I}). \tag{10}$$

When we set $B_t = \mathbf{I}$, we are essentially considering $\mathbf{z}_t$ as an external variable and independent of the U-Net parameters $\theta$, rather than recursively dependent on $\theta$. Otherwise, by the chain rule, it

# Tips: Example

- **Discuss & Exemplify:**

approximated by empirical average when $K$ increases. We note one can also construct a fixed set of $\{u^{(i)}\}_{i\in[N]}$, for example half male and half female to represent a balanced gender distribution. But a fixed split poses a stronger finite-sample alignment objective and neglects the sensitivity of OT.

image $o^{(i)}$ are generated using the same initial noise. We use both CLIP and DINO because CLIP is pretrained with text supervision and DINO is pretrained with image self-supervision. In imple-

- **Provide critical details:**

substantially edited by the DAL. In our implementation, we use the CelebA (Liu et al., 2015) and the FairFace dataset (Karkkainen & Joo, 2021) as external faces. We use the SFNet-20 (Wen et al., 2022) as the face embedding model.

the underlying issue. Moreover, standardizing $A_i$ further stabilizes the optimization process. We note that, to reduce the memory footprint, in all experiments we (*i*) quantize the diffusion model to `float16`, (*ii*) apply gradient checkpointing (Chen et al., 2016), and (*iii*) use DPM-Solver++ (Lu et al., 2022) as the diffusion scheduler, which only requires around 20 steps for T2I generations.

# Tips: Example

- **Discuss & Exemplify**

approximated by empirical average when $K$ increases. We note one can also construct a fixed set of $\{u^{(i)}\}_{i\in[N]}$, for example half male and half female to represent a balanced gender distribution. But a fixed split poses a stronger finite-sample alignment objective and neglects the sensitivity of OT.

image $o^{(i)}$ are generated using the same initial noise. We use both CLIP and DINO because CLIP is pretrained with text supervision and DINO is pretrained with image self-supervision. In imple-

- **Provide critical details:**

substantially edited by the DAL. In our implementation, we use the CelebA (Liu et al., 2015) and the FairFace dataset (Karkkainen & Joo, 2021) as external faces. We use the SFNet-20 (Wen et al., 2022) as the face embedding model.

the underlying issue. Moreover, standardizing $A_i$ further stabilizes the optimization process. We note that, to reduce the memory footprint, in all experiments we (*i*) quantize the diffusion model to float16, (*ii*) apply gradient checkpointing (Chen et al., 2016), and (*iii*) use DPM-Solver++ (Lu et al., 2022) as the diffusion scheduler, which only requires around 20 steps for T2I generations.
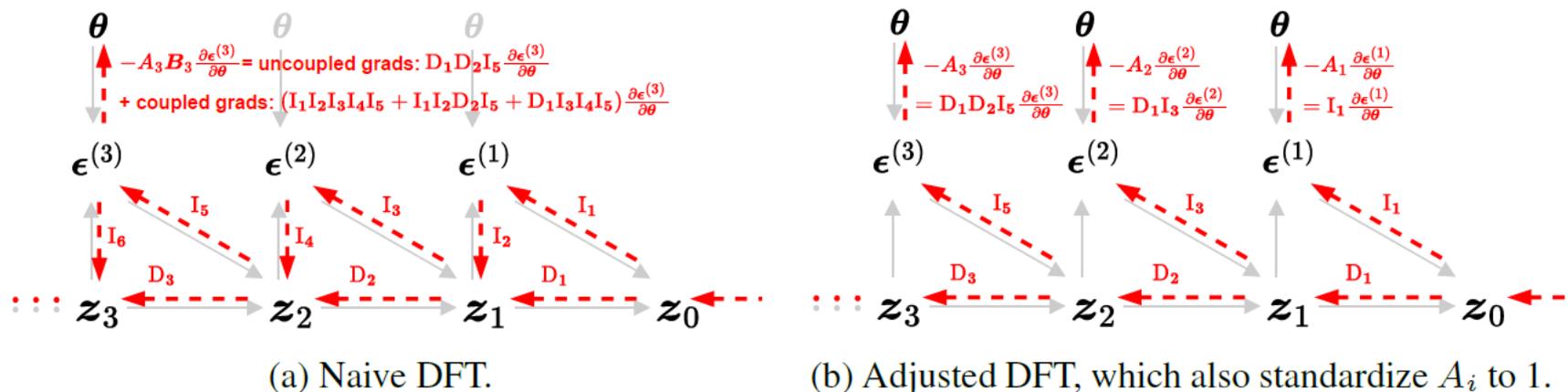
# Tips: Example

- Visualisation



(a) Naive DFT.

(b) Adjusted DFT, which also standardize $A_i$ to 1.

Figure 2: Comparison of naive and adjusted direct finetuning (DFT) of the diffusion model. Gray
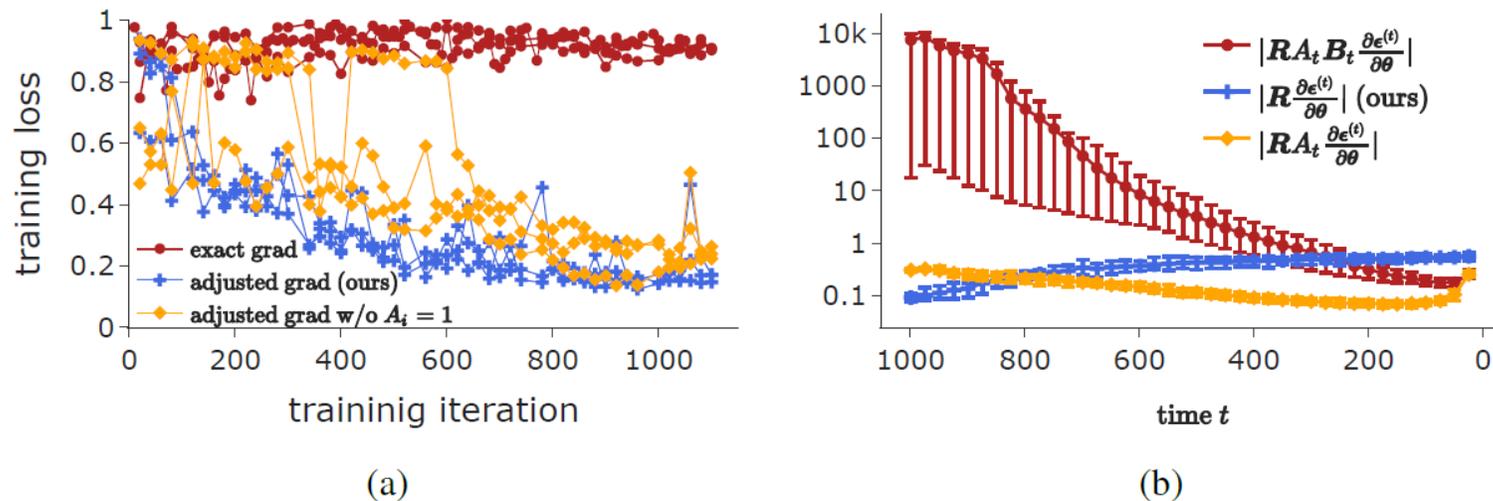


(a)

(b)

Figure 1: The left figure plots the training loss during direct fine-tuning, w/ three distinct gradients.

# Tips: Example

- **Present the algorithm**

- **Release code and data**

---

**Algorithm A.1:** Adjusted DFT of diffusion model

**input** : text encoder $g_\phi$; U-Net $\epsilon_\theta$; image decoder $f_{Dec}$; loss function $\mathcal{L}$; variance schedule $\{\beta_t \in (0,1)\}_{t \in [T]}$ and corresponding $\{\alpha_t\}_{t \in [T]}$, $\{\bar{\alpha}_t\}_{t \in [T]}$; prompt P; diffusion scheduler `scheduler`; inference time step schedule $t_1 = T, t_2, \cdots, t_S = 0$.

```
/* Prepare grad coefficients.                                          */
```
**for** $i = 1, 2, \cdots, S-1$ **do**
$\quad C_i \leftarrow 1/(\frac{1}{\sqrt{\bar{\alpha}_{t_i}}} \frac{\beta_{t_i}}{\sqrt{1-\bar{\alpha}_{t_i}}})$;
**end**

$C \leftarrow [C_1, \cdots, C_{S-1}]/(\prod_{t=1}^{K-1} C_t)^{1/(S-1)}$;
```
/* T2I w/ adjusted gradient                                            */
```
$z_t \leftarrow z_T \sim \mathcal{N}(0,1)$;
**for** $i = 1, 2, \cdots, S-1$ **do**
$\quad t, t_{prev} \leftarrow t_i, t_{i+1}$;
$\quad z'_t \leftarrow \texttt{detach}(z_t)$;
$\quad \epsilon_t \leftarrow f_\theta(g_\phi(\text{P}), z'_t, t_t)$;
$\quad \epsilon'_t \leftarrow \epsilon_t.\texttt{grad\_hook}(g : g \times C[i])$;
$\quad z_{t_{prev}} \leftarrow \texttt{scheduler}(z_t, \epsilon'_t, t, t_{prev})$;
**end**
$x_0 \leftarrow f_{Dec}(z_0)$;

Backpropagate gradient $\frac{d\mathcal{L}(x_0)}{dx_0}$ from generated image $x_0$ to U-Net $\theta$, text encoder $\phi$, or prompt P

---

REPRODUCIBILITY STATEMENT

We provide our source code and various trained fair adaptors for *runwayml/stable-diffusion-v1-5* in `https://github.com/sail-sg/finetune-fair-diffusion`. The pseudocode for

# Tips: Example

- **Proposition and Proof**

**Proposition 1.** *Let $V$ be a set of $n$ nodes with specific degrees $d_1, d_2, ..., d_n$, on which graph $G$ is generated from the configuration model. Let $m = \frac{1}{2}\sum_{i=1}^{n} d_i$ denote the number of edges and $\bar{k} = (\sum_i d_i^2)/(\sum_i d_i)$ be the expected degree when a node is chosen with probability proportional to its degree. As $n \to \infty$, for any node $i \in V$, its local i-quad coefficient satisfies:*

$$\mathbb{E}[I(i)] = \frac{(\bar{k}-1)^2}{2m},$$

*and its local o-quad coefficient satisfies:*

$$\mathbb{E}[O(i)] = \frac{(d_i-1)\cdot(\bar{k}-1)}{2m}.$$

*Proof.* For any open quadriad with node $i$ as an inner node, we denote one outer node by $k$ and another outer node by $l$ (Figure 9a). The probability that this open quadriad is closed equals the probability of having an edge between node $k$ and $l$, which is $(d_k-1)(d_l-1)/2m$ in the configuration

mode. The reason of subtracting 1 from $d_k$ and $d_l$ is that one stub of node $k$ (and node $l$) has already been used in forming the open quadriad.

Now, we show that as $n \to \infty$, $\mathbb{E}[d_k] = \mathbb{E}[d_l] = \bar{k}$. Via stub matching, any node, other than node $i$ and $j$, can form an edge with node $j$ and thus become one outer node of the open quadriad. The probability of node $k$ being this node is proportional to its degree, which is $\frac{d_k}{\sum_{k\in V, k\neq i,j} d_k}$. Therefore, we have $\mathbb{E}[d_k] = \sum_{k\in V, k\neq i,j} d_k \cdot \frac{d_k}{\sum_{k\in V, k\neq i,j} d_k}$. When $n \to \infty$, $\mathbb{E}[d_k] = \sum_{k\in V} d_k \cdot \frac{d_k}{\sum_{k\in V} d_k} = \bar{k}$. Similarly, we have $\mathbb{E}[d_l] = \bar{k}$.

In short, we have:

$$\mathbb{E}[I(i)] = \mathbb{E}[(d_k-1)(d_l-1)/(2m)]$$
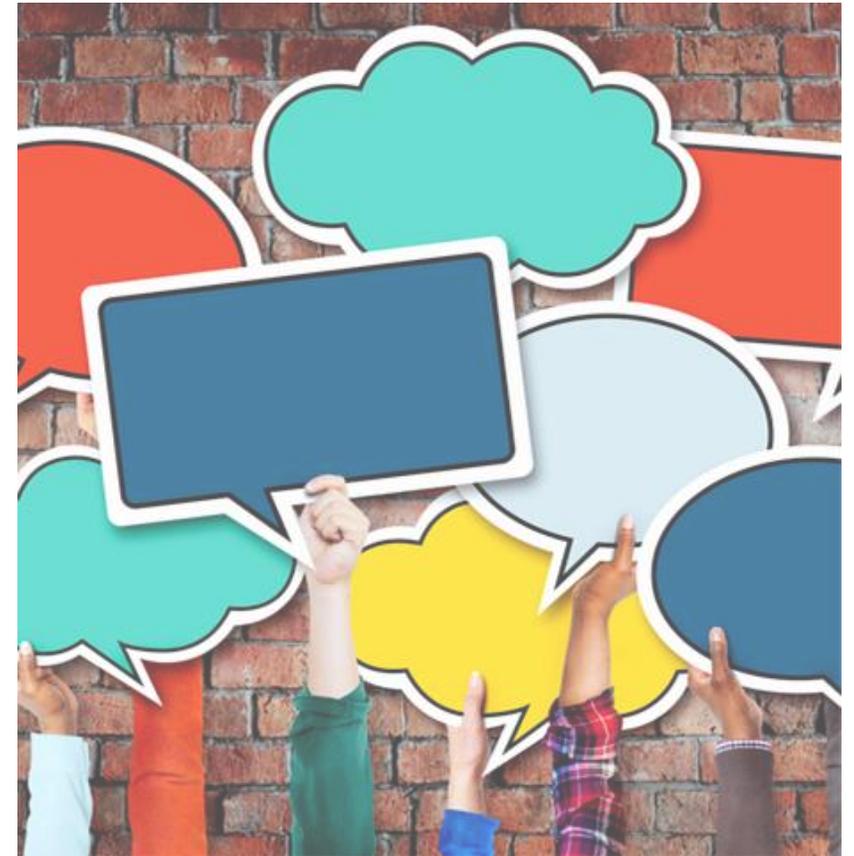$$= \frac{(\mathbb{E}[d_k]-1)\cdot(\mathbb{E}[d_l]-1)}{2m} = \frac{(\bar{k}-1)^2}{2m}.$$

Likewise, for any open quadriad with node $i$ as an outer node, we denote the other outer node by $l$ (Figure 9b). And we have:

$$\mathbb{E}[O(i)] = \mathbb{E}[(d_i-1)(d_l-1)/(2m)]$$
$$= \frac{(d_i-1)\cdot(\mathbb{E}[d_l]-1)}{2m} = \frac{(d_i-1)\cdot(\bar{k}-1)}{2m}.$$

□

# Homework

1. Choose a research paper of your interest, and analyse how its methodology section is written.
   - How does the author justify the design choices?
   - How does the author demonstrate rigour & credibility?
   - Is there a figure included? If yes, how does it enhance the understanding?

2. Sketch your own methodology:
   - Outline the headings and sub-headings.
   - List the reasons for at least two design choices.
   - Create a visual that would appear in the section.

# 4. Publishing Your Research

# Why Publish?

Why Publish

Share findings and contribute to the field

Increase visibility and impact

Build up track record and academic profile

Networking & Collaboration opportunities

Meet degree or work assessment requirements

# Conference & Journal Publications

🎤 **Conference Papers**

Emphasise novelty and timeliness (focus on current research hotspots)

Main purpose is peer communication — presented orally or as posters at the conference

Shorter length

Faster review and publication cycle

Popular in fields like ML and AI

📚 **Journal Papers**

Emphasise theoretical foundation and in-depth, long-term issues

Main purpose is to provide a comprehensive, validated, and permanent record of research

Longer length

Slower review and publication cycle

In most fields, journal papers are considered more prestigious and recognised

# Conference Venues

**International Conferences**
- A* Top-tier venues  [*60/784*]
  - AAAI, ACL, ACMMM, CVPR,  ICML, ICDE, NeurIPS, KDD, VLDB, WWW, etc.
- A: Excellent venues  [*117/784*]
- B: Good venues  [*221/784*]
- C: Acceptable  venues  [*361/784*]

**Australasian Conferences**
- Australasian B: 6/25
  - AJCAI, OZCHI, ACE, ACISP, ADC, ADCS
- Australasian C: 19/25

CORE2023 Summary:
A* - 7.65% of 784 ranked venues
A - 14.92% of 784 ranked venues
B - 28.19% of 784 ranked venues
Australasian B - 0.77% of 784 ranked venues
C - 46.05% of 784 ranked venues
Australasian C - 2.42% of 784 ranked venues
Other - 171 total

Core Ranking

*Total worldwide S&E publication output reached 3.3 million articles in 2022*.
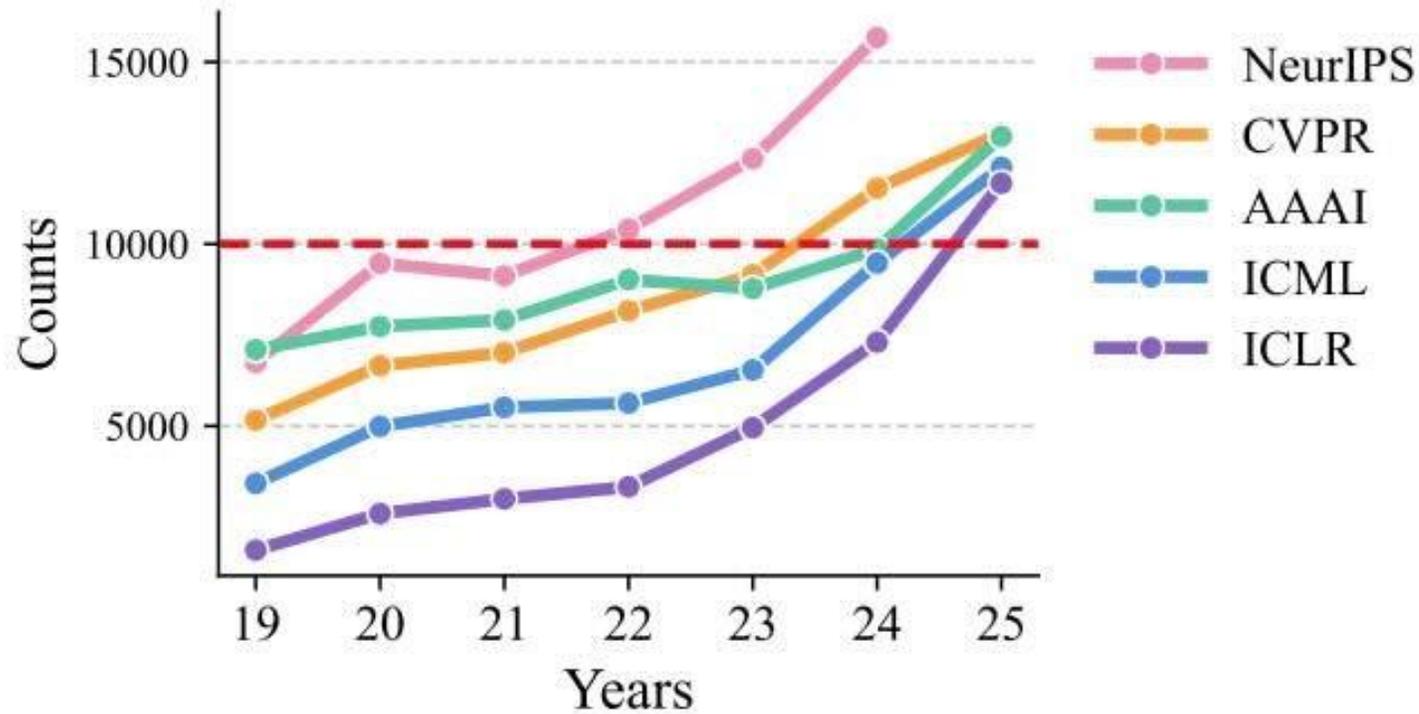
**AI Conference Submissions**

**Figure 1: AI Conference Submission Counts.** The number of paper submissions to most major AI conferences (*e.g.* NeurIPS, CVPR, AAAI, ICML, ICLR) exceeded 10,000 by 2025. For example, there was a 59.8% increase in ICLR submissions in 2025 alone. We forecast similar growth in other venues as well.

source

# Conference Paper Submission and Review Process

[AAAI 2026](#)
- Author Registration and abstract submission (*1 week before submission*)
- Full paper submission (July 12)
- Supplementary material and code (*3 days after submission*)
- Notification of Phase 1 rejections (*7 weeks after submission*)
- Author feedback (*11 weeks after submission*)
- Notification of final decision (*15 weeks after submission*)



The 40th Annual AAAI Conference on Artificial Intelligence
JANUARY 20 – JANUARY 27, 2026 | SINGAPORE

[AJCAI 2025](#)
- Paper submission (June 30)
- Notification (*7 weeks after submission*)



AJCAI
Australasian Joint Conference on Artificial Intelligence 2025
Canberra, Australia
1st Dec - 5th Dec, 2025

# Other Conference Opportunities

## Workshop ([ACMMM 2025](#))

- Shorter and easier to be accepted
- Original research papers, case studies, and position papers

- Personalised Incremental Learning in Medicine
- Workshop on Multimedia Content Analysis in Sports
- Deepfake Forensics: Detection, Attribution, Recognition, and Adversarial Challenges in the Era of AI-Generated Media
- Large Vision Language Model Learning and Applications
- Automotive and Medical Multimedia: Bridging the Gap Between Mobility and Healthcare

## Special Track ([IJCAI 2025](#))

- Focusing on a specific theme or topic
- Peer-reviewed to the same standard as the main track

Main Track
AI, Arts & Creativity
Human-Centred Artificial Intelligence
AI4Tech: AI Enabling Critical Technologies
AI And Social Good
Survey Track
Demonstrations Track

Questions?