

# Assignment - Generative and Collaborative AI

662354\_A25

## Description:

This assignment accounts for 100% of the final mark for this module. You are required to choose one of the two modalities of generative models covered in the module (Text or Image) and complete a series of tasks, as detailed in the following sections. Your work must be presented in a single written report (maximum 3000 words) submitted in PDF or DOCX format. The report should clearly describe the objectives of your work, background research, experimental setup, evaluation methodology, and results. In addition to the report, you must submit your accompanying code. The code will be assessed separately based on the criteria outlined in the task specification table.

**Important:** **Do not** use AI tools to generate or write the content of your report. Doing so may result in a fail or, in the worst, be treated as academic misconduct. If you use AI tools for proofreading only, you must explicitly state this in your report. You must also make sure to properly cite any external sources or tools/models you use in your work.

## Coursework Goal:

In this assignment, your goal is to explore and experiment with one of the two core modalities of generative models introduced in this module: Image Generation or Text Generation. Rather than working on a fixed problem or dataset, you are encouraged to define your own topic of interest, such as *generating photorealistic images from sketches, performing text style transfer, synthesising training data (text or image) for a specific domain, or using large language models for summarisation, translation, or code generation.*

Your task is to test and fine-tune a generative system around your chosen use case, and evaluate its performances and behaviour across three key stages:

1. Baseline Generation using a standard pre-trained model and default prompts
2. Prompt Engineering to refine and control outputs
3. Fine-Tuning the model on a small custom dataset

At each stage, you will evaluate results using appropriate metrics (e.g., FID, CLIP score, BLEU, BERTScore) and reflect on both the quantitative and qualitative differences in output. Alongside this, you will be expected to conduct a short literature review, consider the ethical implications of your work, and submit a clean, modular codebase that reproduces your results.

## Tasks:

The table below outlines the specific tasks that make up this assignment, along with the corresponding marks allocated to each.

Generative and Collaborative AI	
Criteria	Score
<u>Coursework goal:</u> <ul style="list-style-type: none"> <li>Clearly state what you're trying to achieve:           <ul style="list-style-type: none"> <li>Choose one of the two major modalities of generative models taught in class, for example:  <b>Image Generative Models:</b> Generate photorealistic images from sketches, style transfer, generate synthetic data for domain X etc. or,  <b>Large Language Models:</b> Generate product descriptions, Summarize long documents, Generate code from natural language specifications, translate text between languages or styles etc.</li> </ul> </li> <li>Justify why this goal is interesting or relevant.</li> <li>Include a well-defined problem statement and intended outcome.</li> </ul>	<b>5</b>
<u>Literature-Inspired Design:</u> <ul style="list-style-type: none"> <li>Review at least 3 academic or technical papers related to your topic.</li> <li>For each study:           <ul style="list-style-type: none"> <li>Summarize the methodology</li> <li>Discuss the strengths and limitations</li> <li>Mention how it relates to your chosen topic</li> </ul> </li> </ul>	<b>10</b>
<u>Generate Baseline</u> <ul style="list-style-type: none"> <li>Use a standard pre-trained model (e.g., Stable Diffusion 1.5, GPT-2, DistilGPT2 etc.).</li> <li>Choose 3–5 standard prompts.</li> <li>Generate and include the results in the report.</li> </ul>	<b>5</b>
<u>Evaluate Baseline Results</u> <ul style="list-style-type: none"> <li>Use appropriate evaluation metrics, such as:           <ul style="list-style-type: none"> <li>CLIP score</li> <li>FID (Fréchet Inception Distance)</li> <li>IS (Inception Score)</li> <li>BLEU</li> <li>ROGUE</li> <li>BERTScore</li> <li>Human-perceived realism</li> </ul> </li> <li>Report and interpret the scores.</li> </ul>	<b>5</b>
<u>Prompt Engineering</u> <ul style="list-style-type: none"> <li>Modify your original prompts using prompt engineering techniques.</li> <li>Generate a new set of outputs using the same model.</li> <li>Discuss how prompt changes influenced the outputs.</li> </ul>	<b>15</b>

<u>Evaluate Engineered Prompts</u>	<ul style="list-style-type: none"> <li>• Use the same evaluation metric(s) from Criteria 4.</li> <li>• Compare with the baseline scores and provide a short analysis.</li> </ul>	<b>5</b>
<u>Fine-Tune the Generator (e.g., using LoRA, RAG etc.)</u>	<ul style="list-style-type: none"> <li>• Fine-tune the model using:           <ul style="list-style-type: none"> <li>- A small custom dataset</li> <li>- Techniques like LoRA, RAG etc.</li> </ul> </li> <li>• Document:           <ul style="list-style-type: none"> <li>- The dataset used</li> <li>- Fine-tuning procedure</li> <li>- Training duration</li> <li>- Parameters or hyperparameters</li> </ul> </li> </ul>	<b>15</b>
<u>Generate and Evaluate Fine-Tuned Outputs</u>	<ul style="list-style-type: none"> <li>• Evaluate using the same metrics.</li> <li>• Include generated results and scores.</li> <li>• [Optional]: Use the same original and engineered prompts on the fine-tuned model.</li> </ul>	<b>10</b>
<u>Compare &amp; Contrast Results</u>	<ul style="list-style-type: none"> <li>• Present a side-by-side comparison of the results:           <ul style="list-style-type: none"> <li>- Pretrained vs. Prompt-Engineered vs. Fine-Tuned</li> </ul> </li> <li>• Provide both:           <ul style="list-style-type: none"> <li>- Quantitative analysis (scores)</li> <li>- Qualitative analysis (visual or subjective impressions)</li> </ul> </li> <li>• Reflect on what worked, what didn't, and why.</li> </ul>	<b>15</b>
<u>Ethical Considerations:</u>	<ul style="list-style-type: none"> <li>• Use this section to discuss the ethical issues that might be involved in the chosen area of research or/and the dataset used. For example:           <ul style="list-style-type: none"> <li>- Are there risks of bias, or misinformation in the model's output?</li> <li>- Does the dataset contain sensitive, private, or potentially harmful content?</li> <li>- Could the generated content be misused (e.g., deepfakes, fake news, impersonation etc.)?</li> <li>- Is there a risk of reinforcing harmful stereotypes, language, or social biases?</li> </ul> </li> </ul>	<b>5</b>
<u>Code Submission</u>	<ul style="list-style-type: none"> <li>• Submit a clean, well-commented Jupyter Notebook.</li> <li>• Code should be modular and reproducible.</li> <li>• Outputs should be matching with the report submitted.</li> </ul>	<b>10</b>



Task	Criteria	Max Points	First	2:1	2:2	Third	Poor
1	Goal	5	Goal is precise & measurable. Modality fit is justified; clear success criteria & scope, along with assumptions/risks, are stated. <b>5 pts</b>	Clear goal & modality fit; success criteria partly quantified; minor scope or risk gaps. <b>4 pts</b>	General goal (e.g., “improve quality”); modality fit implied; outcomes loosely defined. <b>3 pts</b>	Vague/over-broad; success criteria missing; weak link to modality/problem. <b>2 pts</b>	No coherent goal or off topic. <b>0–1 pt</b>
2	Inspiration from Literature	10	≥3 strong sources; each has a 1–2 paragraph critique covering method, data, compute, strengths/limits, an explicit “design takeaway” mapped to your pipeline. <b>10 pts</b>	3 sources summarised with some critique & takeaways; a few method/data gaps. <b>8 pts</b>	2–3 sources; mostly descriptive; weak or generic takeaways. <b>6 pts</b>	Superficial summaries; unclear relevance; missing key details. <b>4 pts</b>	Missing/irrelevant or uncited. <b>0–2 pts</b>
3	Generate Baseline	5	Reproducible run of a standard model with 3–5 well-rationalised prompts; show seeds, version, inference params; present 6 – 8 representative outputs. <b>5 pts</b>	Baseline runs; prompts listed; minor omissions (e.g., missing seed or version). <b>4 pts</b>	Minimal run; prompts not justified; few/low-quality exemplars. <b>3 pts</b>	Unclear config; patchy results; not obviously baseline. <b>2 pts</b>	Not run/failed without evidence. <b>0 – 1 pt</b>
4	Evaluate Baseline Results	5	Correct metrics & rationale, with formulas or tool refs; table of scores + 3–5 sentence interpretation per metric (what improves what, limits). <b>5 pts</b>	Appropriate metrics computed; brief interpretation; minor correctness/justification gaps. <b>4 pts</b>	Some metrics or only human judgement; surface-level discussion. <b>3 pts</b>	Misapplied, incomparable, or not reproducible metrics. <b>2 pts</b>	No evaluation. <b>0–1 pt</b>

5	Prompt Engineering	15	Systematic interventions (e.g., role/context framing, constraints, few-shot, ordering, self-critique, temperature/top-p changes). Clear hypotheses; before/after grids; prompt templates versioned; ablation (which edit yielded which change). <b>15 pts</b>	Several sensible edits; persuasive before/after; some hypothesis testing. <b>12 pts</b>	Ad-hoc edits; qualitative comments; limited evidence. <b>9 pts</b>	Minimal edits; weak causal link to changes. <b>6 pts</b>	Not attempted or non-evidenced. <b>0-3 pts</b>
6	Evaluate Engineered Prompts	5	Same metrics reused; neat table: Baseline vs Engineered; effect sizes or % change; short analysis of trade-offs (e.g., diversity↑, faithfulness↓). <b>5 pts</b>	Table/plot comparing scores; surface-level discussion. <b>4 pts</b>	Basic comparison (numbers only), thin insight. <b>3 pts</b>	Partial/inconsistent comparison. <b>1-2 pts</b>	No comparison. <b>0 Pts</b>
7	Fine-Tune the Generator	15	Clear FT strategy. Dataset stated with; train/val split; batch size, steps/epochs, lr/scheduler, reg, compute/time; seeding; why these choices; failure modes & mitigations. <b>15 pts</b>	Sensible FT; adequate hyper-params and data notes; limited justification. <b>12 pts</b>	Minimal FT (few steps or unclear method); missing key params. <b>9 pts</b>	Execution issues; unjustified/opaque setup. <b>6 pts</b>	Not attempted. <b>0-1 pts</b>
8	Generate and Evaluate Fine-Tuned Outputs	10	Re-run same metrics; show qualitative samples; if applicable, reuse baseline/engineered	Metrics + samples; moderate interpretation. <b>8 pts</b>	Results present, limited or uneven analysis. <b>6 pts</b>	Weak link from FT to outcomes; little evidence. <b>4 pts</b>	Missing results/eval. <b>0-2 pts</b>

			prompts; discuss generalisation vs overfit, and cost/latency impacts. <b>10 pts</b>				
9	Compare & Contrast Results	15	One consolidated side-by-side table/figure for Pretrained vs Prompt-Engineered vs Fine-Tuned; highlight wins/losses, variance/Confidence Interval if relevant; integrate qualitative and quantitative; reflective “what worked/why/limits/next” . <b>15 pts</b>	Good, consolidated comparison; mostly quantitative with some reflection. <b>12 pts</b>	Basic tables/plots; limited analysis and discussion. <b>9 pts</b>	Fragmented or narrative-only comparison. <b>6 pts</b>	No coherent synthesis. <b>0–3 pts</b>
10	Ethical Considerations	5	Concrete risks (bias, stereotyping, misinformation, IP/privacy, misuse e.g., deepfakes) tied to your data/model; feasible mitigations (filters, safety prompts, consent, watermarking, eval); license/attribution; limitations. <b>5 pts</b>	Good coverage with minor omissions or generic mitigations. <b>4 pts</b>	General statements; few specifics or missing mitigations. <b>3 pts</b>	Some mentions; no concrete actions or mitigations. <b>1–2 pts</b>	Absent. <b>0 pt</b>
11	Code Submission	10	Single clean Jupyter: modular functions, docstrings, config cell; fixed seeds; clear Run-	Mostly reproducible; minor manual steps. <b>8 pts</b>	Runs with tweaks, partial mismatch with report. <b>6 pts</b>	Hard to run; unclear structure; missing deps. <b>4 pts</b>	Not provided/non-functional. <b>0–2 pts</b>



			All; notebook re-generates report figures/tables; saved artifacts. <b>10 pts</b>				
--	--	--	--	--	--	--	--