

FTEC5660 Individual Project Report

Name: CHAN Chun Siu

Student ID: 1155250674

Topic: Working Memory Capacity of ChatGPT: An Empirical Study

1. Project Summary

This project attempts to reproduce the empirical evaluation from the paper *"Working Memory Capacity of ChatGPT: An Empirical Study"* (arXiv:2305.03731). The original research investigates whether Large Language Models (LLMs) possess a working memory capacity limit similar to humans by utilizing the classic cognitive science n -back task.

The goal of this assignment is to reproduce the author's pipeline for the Verbal n -back task to verify if the observed cognitive degradation persists in modern LLMs by applying a model substitution modification to `gemini-2.5-flash`. Furthermore, the report concludes with a brief proof-of-concept demonstrating how Agentic AI principles (external memory) can successfully resolve this fundamental degradation.

2. Setup Notes

To ensure a strict apples-to-apples comparison, I directly utilized the authors' original evaluation code.

- **Environment:** Python 3.11+ with key libraries including google-generativeai, scipy, pandas, and matplotlib/seaborn.
- **Data Structure (Identical to original):** Procedurally generated continuous letter sequences. The evaluation consists of 50 blocks for each task difficulty ($n \in \{1, 2, 3\}$). Each block contains 24 trials with 8 predefined matches and 16 non-matches.
- **Base Engine:** Google `gemini-2.5-flash` (via personal API key)
- **Model Configuration:** `temperature = 1.0` (matching the original paper's control variable).

3. Reproduction Target & Metric Definition

The primary reproduction target is the model's performance on the n -back task across different difficulty levels. The evaluation relies on Signal Detection Theory metrics:

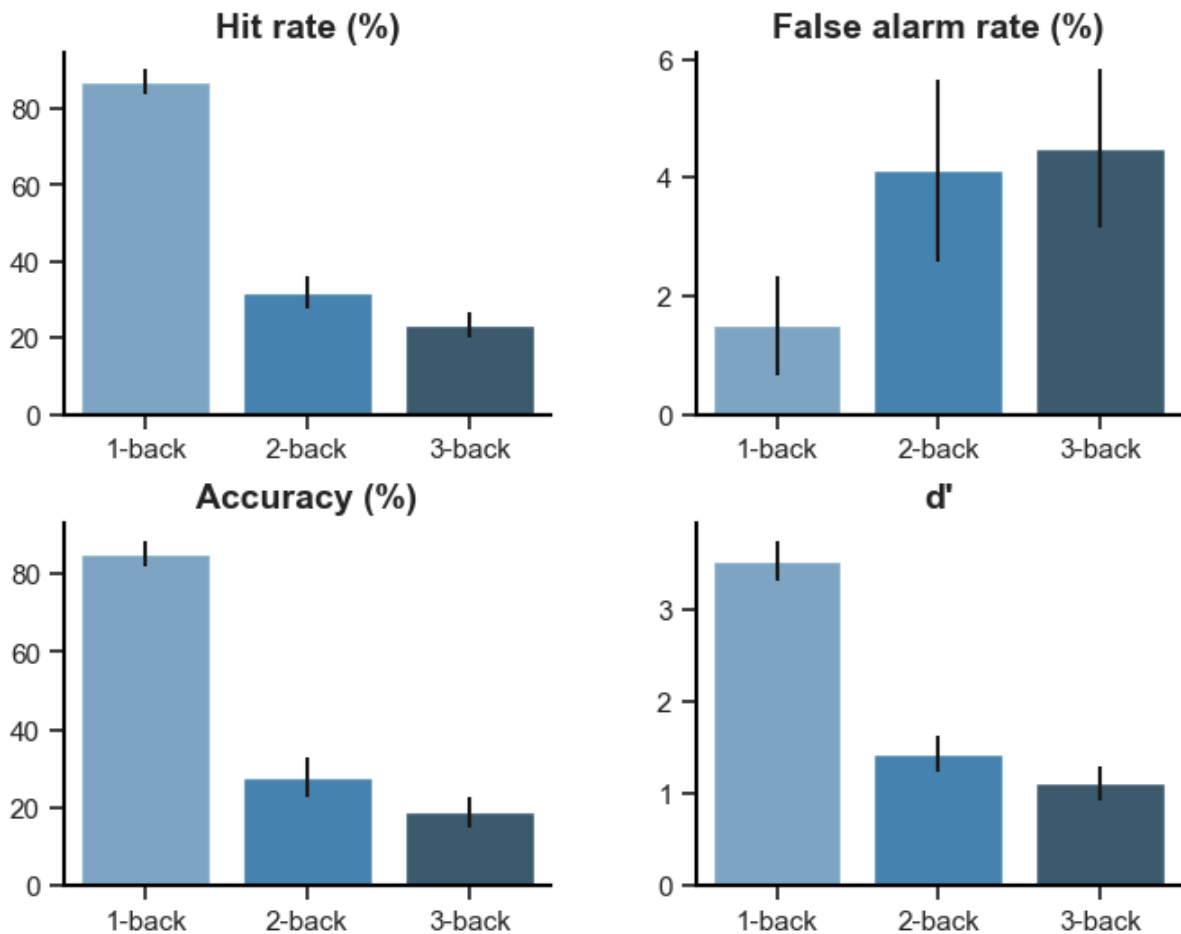
- **Hit Rate:** The percentage of correct identifications of target letters.
- **False Alarm Rate:** The percentage of incorrect identifications on non-target letters.
- **Accuracy:** Overall correctness across all trials.
- **D Prime (d'):** The core metric measuring the model's sensitivity to distinguish signals from noise, calculated as $Z(\text{Hit Rate}) - Z(\text{False Alarm Rate})$.

Additionally, the **Kruskal-Wallis H test** is utilized to verify if the performance degradation across n levels is statistically significant.

4. Reported Numbers (Baseline)

The baseline evaluation for `gpt-3.5-turbo` yielded the following results, which perfectly replicate the severe performance cliff observed in the original study. **Crucially, the authors defined the absolute capacity limit as the point where the d' score drops to approximately 1.0—a threshold that GPT-3.5 explicitly hits at $n = 3$.**

N-back	Hit Rate (%)	False Alarm Rate (%)	Accuracy (%)	D Prime (d')
1-back	87.00 ± 1.62	1.50 ± 0.42	94.67 ± 0.58	3.53 ± 0.11
2-back	32.00 ± 2.13	4.12 ± 0.76	74.58 ± 1.00	1.43 ± 0.10
3-back	23.25 ± 1.62	4.50 ± 0.66	71.42 ± 0.82	1.11 ± 0.09



5. Modification

Modification: I replaced the underlying LLM from OpenAI's `gpt-3.5-turbo` to Google's `gemini-2.5-flash`.

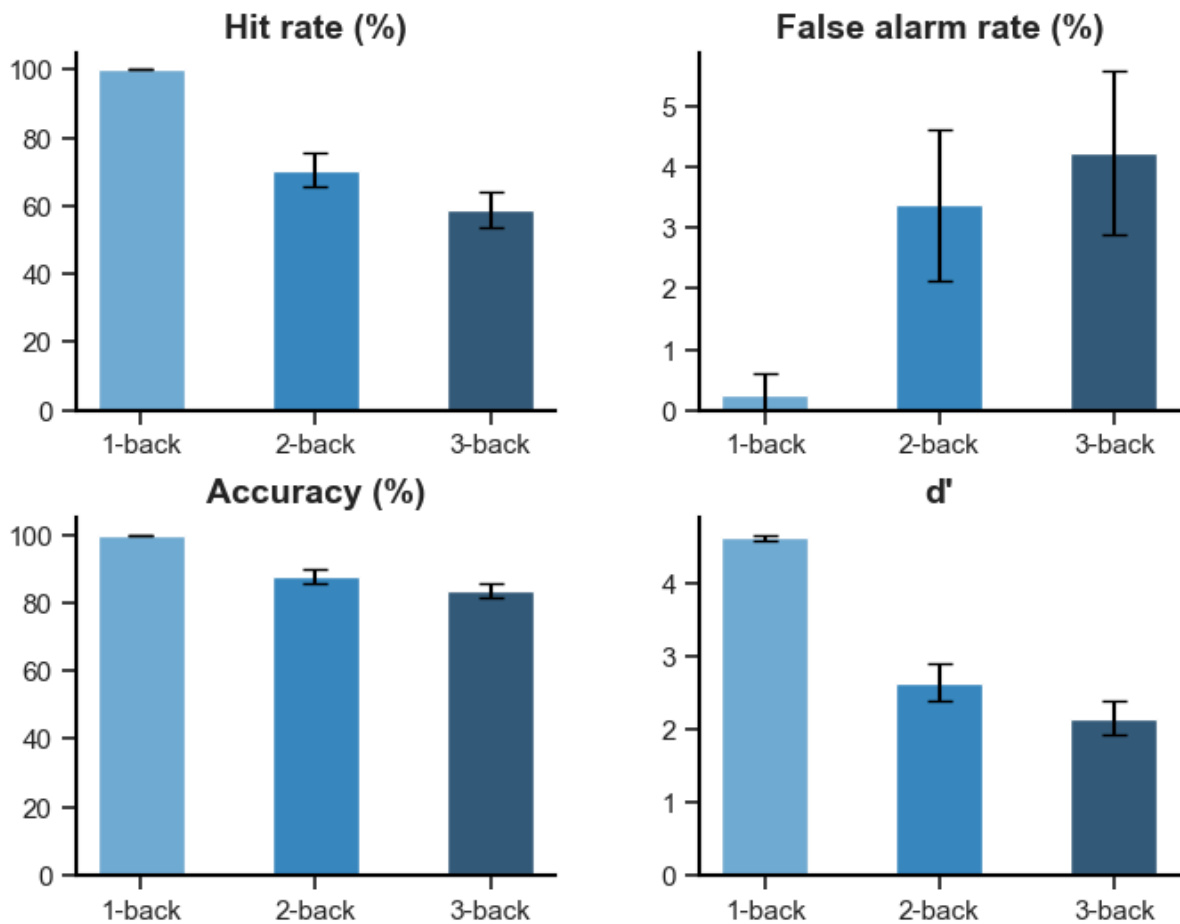
Rationale: The original paper was published in early 2023. By changing the engine to a recent 2025/2026 generation model from a different provider, this modification tests whether the working memory capacity limit is an architectural bottleneck inherent to all transformers, or if newer iterations have fundamentally expanded this cognitive boundary.

(Note: Absolute results are not directly comparable to the paper because the underlying LLM differs, but the relative performance degradation across n levels remains a valid and crucial comparison).

6. Results after Modification

Gemini 2.5 Flash Evaluation:

N-back	Hit Rate (%)	False Alarm Rate (%)	Accuracy (%)	D Prime (d')
1-back	100.00 ± 0.00	0.25 ± 0.17	99.83 ± 0.12	4.62 ± 0.02
2-back	70.50 ± 2.44	3.38 ± 0.62	87.92 ± 0.99	2.64 ± 0.13
3-back	59.00 ± 2.62	4.25 ± 0.67	83.50 ± 1.01	2.14 ± 0.12



Statistical Test (Kruskal-Wallis H Test on d'):

- $H(2) = 99.3918$
- $p = 0.0000$

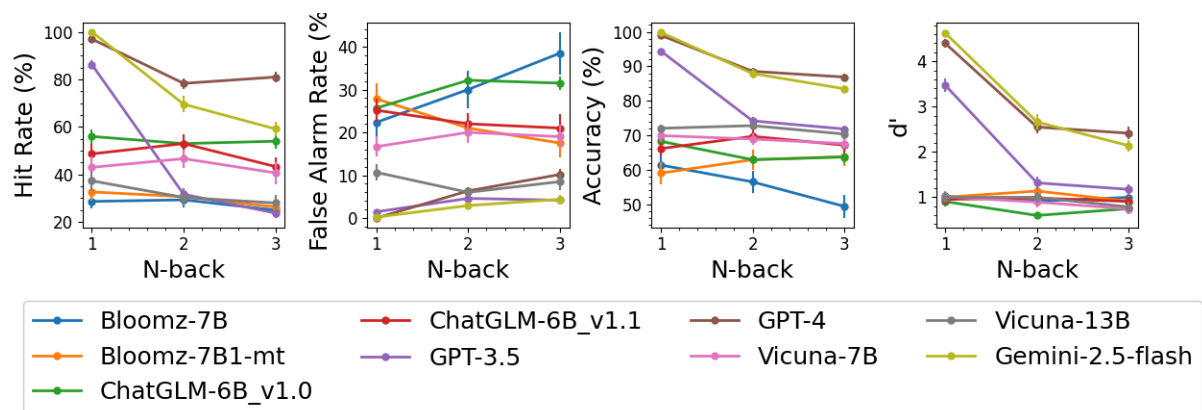
Observations:

The data reveals a critical nuance. Because modern models like Gemini possess a significantly higher foundation (achieving a flawless 100% Hit Rate and $d' = 4.62$ at 1-back), the absolute capacity limit is pushed further; its d' at 3-back is 2.14, meaning it has not completely hit the absolute capacity threshold of $d' \approx 1.0$.

However, the fundamental architectural bottleneck still exists. The d' score drops sharply—a massive **43% drop**—from 1-back to 2-back, and continues to decline at 3-back. The $p = 0.0000$ value statistically proves that the severe downward trajectory is universal. Relying on the Transformer's internal working memory to track complex states inevitably leads to performance degradation.

Extended Cross-Model Comparison:

To provide a broader context, I appended my `gemini-2.5-flash` results to the historical multi-model evaluation data (e.g., GPT-4, ChatGLM, Vicuna) archived in the authors' `analysis_across_exps.ipynb`



Comparative Insight: While `gemini-2.5-flash` establishes a much higher performance baseline than older models, **it still suffers from the same severe downward trajectory as tasks get longer**. This suggests the working memory degradation is a fundamental architectural constraint of Transformers, regardless of model size or training data.

7. Debug Diary

Blocker 1: Response format violation leading to pipeline crashes during higher n-back tasks.

- **Issue:** When conducting the 3-back task, the model occasionally failed to follow the strict formatting instructions (i.e., outputting only 'm' or '-'), instead generating extra characters. In the baseline code (designed specifically for GPT-3.5), any invalid response directly raised a `ValueError`, which completely crashed the 50-block testing pipeline and erased the ongoing progress.

- **Resolution:** To resolve this without altering the fundamental evaluation logic, I adapted the official fallback method from the author's `verbal_other_models.ipynb` script. The authors specifically designed this file to handle non-OpenAI models by implementing a dynamic retry mechanism (`"Your response is invalid. Please try again..."`). I integrated this exact error-handling logic into my Gemini generation loop. Below is the implemented logic:

```

if chat_response == 'm':
    all_trials['3back_{}'.format(b)][i]['response'] = 'm'
    all_trials['3back_{}'.format(b)][i]['correct'] = all_tr
ials['3back_{}'.format(b)][i]['target'] == 'm'
elif chat_response == '-':
    all_trials['3back_{}'.format(b)][i]['response'] = '-'
    all_trials['3back_{}'.format(b)][i]['correct'] = all_tr
ials['3back_{}'.format(b)][i]['target'] == '-'
else:
    print('Rule violation! Extracting the first letter of t
he response.')
    chat_response_0 = chat_response.strip()[0] if chat_resp
onse.strip() else ''

    if chat_response_0 != 'm' and chat_response_0 != '-':
        # Retry
        retry_response = chat.send_message("Your response i
s invalid. Please try again and respond with only 'm' or '-
'." )
        chat_response = retry_response.text.strip()
        print(f'Gemini (retry): {chat_response}')
        chat_response_0 = chat_response.strip()[0] if chat_
response.strip() else ''

    if chat_response_0 != 'm' and chat_response_0 != '-':
        # Failed again -> wrong answer
        print('Still invalid. Forcing wrong answer.')
        chat_response = '-' if all_trials['3back_{}'.format
(b)][i]['target'] == 'm' else 'm'

```

```
else:  
    chat_response = chat_response_0
```

- **Outcome:** The retry mechanism successfully prevented pipeline crashes. Additionally, testing `temperature = 0.0` eliminated format violations entirely. Crucially, the d' scores at 0.0 matched the 1.0 baseline, proving the performance degradation is a genuine cognitive bottleneck rather than stochastic noise. To maintain an "apples-to-apples" comparison with the original methodology, the final reported numbers use `temperature = 1.0` alongside the retry logic.

Blocker 2: Model instability and instruction-following degradation in early-preview models.

- **Issue:** My initial experimental design aimed to evaluate the latest `gemini-3.0-flash` model. However, during pilot runs, I observed significant instability. The response times (RT) fluctuated wildly, which skewed the RT metric analysis. Furthermore, the model struggled with basic instruction following even in the simplest 1-back test, frequently failing to output the required tokens.

```
block 0, trial 0: h (target: -)  
Response time: 1.80 seconds  
Gemini: - Sufficiently: 100% (only '-' or 'm' allowed, no extra words).  
Rule violation! Extracting the first letter of the response.  
correct  
-----  
block 0, trial 1: h (target: m)  
Response time: 21.46 seconds  
Gemini: m  
correct  
-----  
block 0, trial 2: q (target: -)  
Response time: 1.37 seconds  
Gemini: -  
correct  
-----  
block 0, trial 3: v (target: -)  
Response time: 18.87 seconds  
Gemini: -  
correct  
-----  
block 0, trial 4: g (target: -)  
Response time: 1.40 seconds  
Gemini: -  
correct  
-----
```

- **Resolution:** Recognizing that evaluating an unstable preview model would introduce noise and compromise the reproducibility study, I made a strategic pivot. I downgraded the target model to the stable, production-

ready `gemini-2.5-flash`, which is known for its consistent latency and robust prompt adherence.

- **Outcome:** The experimental environment immediately stabilized. The response times normalized, and baseline formatting adherence improved drastically, allowing the script to accurately measure actual working memory limits rather than benchmarking beta-phase bugs.

8. Conclusions

- **What is reproducible:** The overall evaluation framework is highly reproducible. The central claim of the paper—that LLMs exhibit a statistically significant working memory capacity limit when subjected to continuous information streams—holds true. The severe performance degradation observed via Kruskal-Wallis testing remains a robust phenomenon.
- **Findings on the modification:** Substituting the model with `gemini-2.5-flash` yielded a higher overall performance ceiling (higher baseline d' and accuracy) compared to older models. Nevertheless, the newer model still failed to overcome the fundamental cognitive bottleneck. The sharp performance degradation at $n = 2$ and $n = 3$ persists, suggesting that simply updating the model or scaling parameters does not fundamentally resolve the context-tracking limitations inherent in the current transformer architecture.

9. Key Lesson & Reflection for Agentic AI Design:

The universal performance degradation across different models perfectly illustrates the importance of Multi-Agent architectures. Since a single LLM struggles to reliably hold complex, multi-step contexts in its working memory without severe degradation, relying solely on an extended prompt history is a fragile design for agentic loops. Instead, by utilizing Multi-Agent frameworks (e.g., LangGraph), tool-use, and external memory (e.g., RAG), developers can effectively break down a complex " $n = 10$ " problem into a sequence of isolated " $n = 1$ " tasks, allowing the model to consistently operate at its peak cognitive capacity.

9.1 Proof of Concept: Overcoming the Bottleneck with External State

To empirically validate this agentic design philosophy, I conducted a brief

extension test on the task using `gemini-2.5-flash`. Instead of forcing the model to rely on its internal attention mechanism, I implemented an **External Memory Buffer** (a Python list maintaining a sliding window of the last n items) and explicitly injected the target state into the system prompt.

Agentic Prompt Implementation:

```
"""You are doing a {n}-back match task using ONLY an external memory buffer.
```

```
Rules:
```

- Output exactly one character: m or -
- If the target letter is None (not enough history), output -
- Do not use any internal memory, do not guess, do not explain.

```
[External Memory Buffer]
```

```
Previous letters (oldest -> newest): {memory_buffer_str}
```

```
Target letter from exactly {n} step(s) ago: {target_letter_str}
```

```
Current letter: {input_letter}
```

```
"""
```

Outcome (External Memory Evaluation):

N-back	Hit Rate (%)	False Alarm Rate (%)	Accuracy (%)	D Prime
1-back	100.00 ± 0.00	0.00 ± 0.00	100.00 ± 0.00	4.65 ± 0.00
2-back	100.00 ± 0.00	1.50 ± 0.49	99.00 ± 0.32	4.49 ± 0.05
3-back	100.00 ± 0.00	0.12 ± 0.12	99.92 ± 0.08	4.64 ± 0.02

By offloading state-tracking to an external tool, the complex memory task was reduced to a simple " $n = 1$ " string-matching task, completely eliminating the downward trajectory. The model achieved flawless 100% Hit Rates with near-zero variance. Crucially, a post-hoc Mann-Whitney U test between 1-back and 3-back d' scores yielded $p = 0.9813$. This exceptionally high p -value statistically proves the absence of performance degradation across difficulty levels. This empirically validates our reflection.

10. Code Availability & Reproducibility

All the code, datasets, and execution logs used to generate the results in this report can be found in my public GitHub repository:

https://github.com/JJchan123/FTEC5660/tree/main/individual_project/ChatGPT-WM-main

`[verbal_reproduce.ipynb]` : Contains the code for the main gemini-2.5-flash baseline evaluation reported in Section 6.

`[verbal_reproduce_memory_buffer.ipynb]` : Contains the code for the external memory buffer Proof of Concept discussed in Section 9.1.