# Agenda

- **Topic Introduction**
- **Background**
- **Goal**
- **Evaluation Criteria**
- **Dataset Description**
- **Notice**
- **Submit Results**
- **Q&A**

2025 TSMC IT
CareerHack

# Speaker



徐梓軒

EPCD / Enterprise and Public Cloud Department

**IT 工程師**
education
**交大AI所**
experience
**台積電 IT 2 年**
- EPCD 工程師

**thhsuy@tsmc.com**
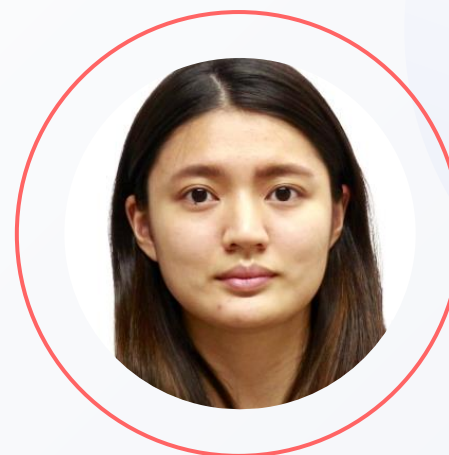
# Meet Our Mentors

**蘇予宣**
Chris Su
IT Engineer

企業及公有雲建構部
EPCD
**yxsuc@tsmc.com**

**陳德優**
Luis Chen
Technical Manager

企業及公有雲建構部
EPCD
**tychenx@tsmc.com**

**李毓簫**
Daisy Li
IT Engineer

企業及公有雲建構部
EPCD
**yhlizzi@tsmc.com**

**陳伯維**
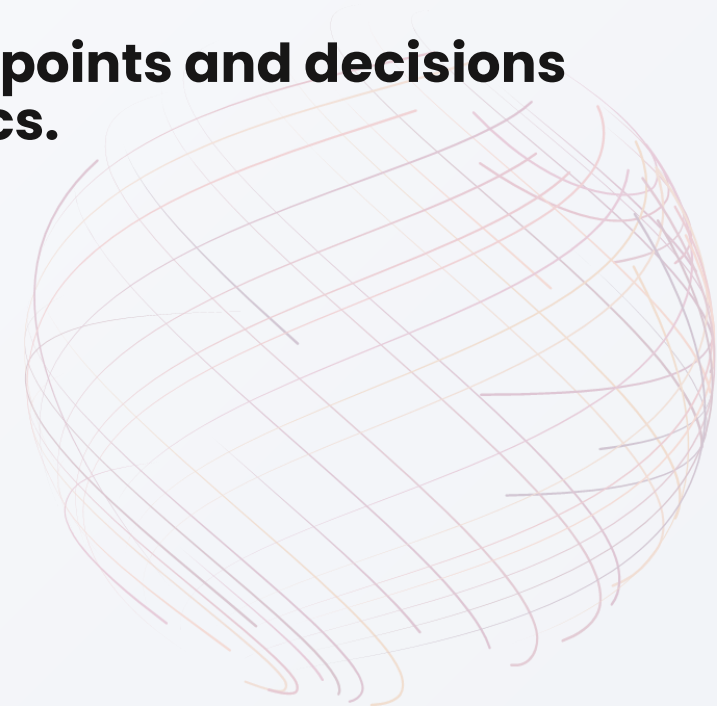Martin Chen
IT Engineer

企業及公有雲建構部
EPCD
**bwchent@tsmc.com**

# Background

- **As new factories expand globally, <span style="color:red">cross-language discussions</span> are often necessary for internal meetings, handovers, and client meetings.**

- **The abundance of <span style="color:red">proper nouns</span> in the company hinders from understanding discussions and key points.**

- **Without <span style="color:red">meeting minutes</span>, it is easy to forget the key points and decisions that have been discussed when there are many topics.**

Security C - TSMC Secret

# Three major system requirements

**Real-Time Speech To Text Translation**

ASR, Text Translation

**Company-Specific Terminology Detector**
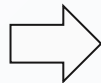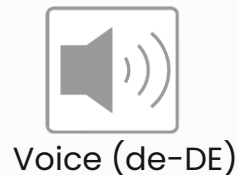
Key Word Detect, Specific Translate

**Clerk Of Meetings**

Search, Summarize

Security C - TSMC Secret

# Task 1: Real-Time Speech To Text Translation

- **Objective: The system should provide instantaneous <span style="color:red">speech to text translation</span> between multiple languages to facilitate seamless communication.**

- **Evaluation:**
  - Ability to handle multiple languages (Chinese, English, Japanese, German)
  - Obtain correctly translation text result
  - Real-time speech to text translation

Voice (de-DE) →

en-US: I worked the night shift yesterday and handed over the matter to LISA.
zh-TW: 我昨天值大夜班，有把事情交接給LISA了
ja-JA:私は昨日夜勤をしていたので、その件をLISAに引き継ぎました。
de-DE: ich hatte gestern Nachtschicht und habe die Angelegenheiten an Lisa übergeben

# Task 2: Company-Specific Terminology Detector

- **Objective: Integrate GenAI capabilities to understand and appropriately translate TSMC-specific jargon and technical terms.**

- **Evaluation:**

  - Detect technical terms and translate them accurately
  - Obtain correctly translation text result
  - Attached the description of professional terminology

我昨天值大夜班 ┊ I worked the night shift yesterday ✔
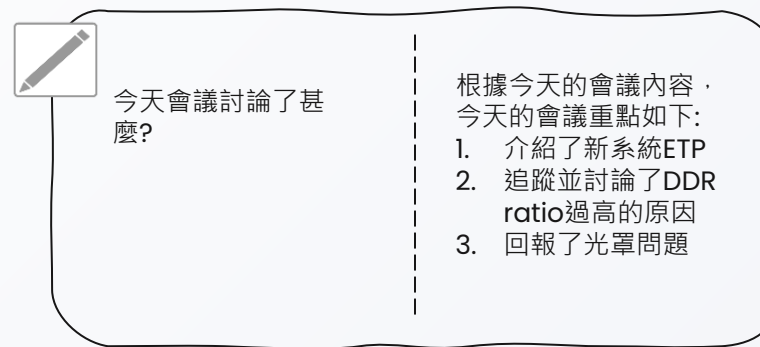
I worked the big night shift yesterday ✘

我們新的系統ETP主要是針對跨國廠區的會議及交接時使用 ┊ ETP: Enterprise Translation Platform，翻譯整合平台，提供安全且即時的翻譯服務，包含會議紀錄及查詢等功能。

20
25 | TSMC IT **CareerHack**
Disruptive Innovation in AI Era

# Task 3: Clerk Of Meetings

- **Objective: Ability to record and transcribe meetings for future reference , system should enable users to search and summarize through the transcribed content after meeting.**

- **Evaluation:**
  - UI for retrieving historical transcriptions or performing transcription searches
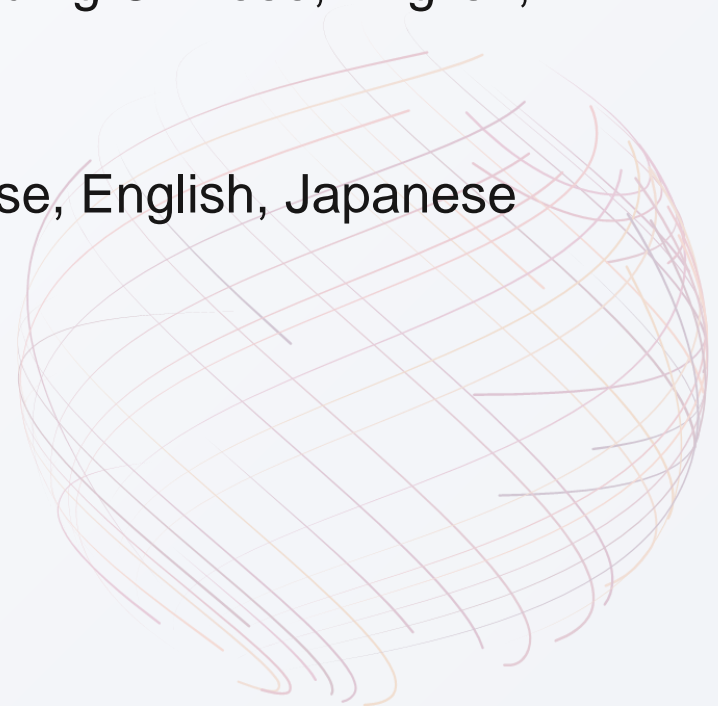  - Ease of use, including the clarity of the interface and search functionality

今天會議討論了甚麼?

根據今天的會議內容，今天的會議重點如下:
1. 介紹了新系統ETP
2. 追蹤並討論了DDR ratio過高的原因
3. 回報了光罩問題

# Evaluation Criteria

| Item | Evaluation Description | Usage Dataset | Score |
|------|------------------------|---------------|-------|
| Real-Time Speech To Text Translation | • Real-time translation function, can select specify language.<br>• Output your Transcript, Chinese is the main language for scoring. | Testing dataset | 10 |
| Company-Specific Terminology Detector | • Detect and accurately translate proper nouns based on the provided dataset, 10~15 proper nouns includes in audio file.<br>• List the proper nouns you detected.<br>• Display the explanations of proper nouns on front-end. | Testing dataset | 15 |
| Clerk Of Meetings | • Store and organize meeting minutes, display with user input language.<br>• The presentation results need to include the following items: meeting summary, assigned tasks. | Training dataset<br>Testing dataset | 15 |
| System Architecture | • Including completion, design flow, UI or full system. | - | 25 |
| Creativity | • In addition to data processing and analysis, what other usage or scenario can be used? | - | 10 |
| Presentation | • The solution should be presented in a clear and concise manner, with appropriate documentation. | - | 10 |
| Live Demo | • Live demo your solution with demo audio, the audio file lasts about 1 minute.<br>• Output your Transcript, Chinese is the main language for scoring. | Demo dataset | 15 |

Pie chart:
- 10% Real-Time Speech To Text Translation
- 15% Company-Specific Terminology Detector
- 15% Clerk Of Meetings
- 25% System Architecture
- 10% Creativity
- 10% Presentation
- 15% Live Demo

# Dataset Description

- **Training.zip:**
  - Training.wav: Recording during the meeting, including Chinese, English, Japanese and German.
  - Knowledge Dataset.xlsx: Special terminology data set, including Chinese, English, Japanese and German.

- **Testing.zip:**
  - Testing.wav: Recording during the meeting, including Chinese, English, Japanese and German.

# Knowledge Dataset.xlsx

## Contains proper nouns and corresponding translations in four languages

Sheet: cmn-Hant-TW, en-US, ja-JP, de-DE



## Contains Transcript and Chinese translation of Training.wav, Red is type fab, Blue is type IT and GCP

Sheet: Training wav

12

# Notice

- **Remember turning off GCP environment once you don't perform anything.**

- **Download dataset from "careerhack2025-icsd-resource-bucket" GCP bucket with gsutil command.**
  - gsutil cp gs://careerhack2025-icsd-resource-bucket/Training.zip gs://{your bucket name}/

- **Resource:**
  - Cloud Storage & Artifact Registry: Each team has 200 GB to store your data.
  - Compute Engine, Cloud Run, Cloud SQL, GKE: Each team has 60 hours to use.
  - 5 points will be deducted if GPU(GKE) exceeds one hour

- **Teamwork is important, which is also important in the workplace.**

- **Please send your question to careerhack@tsmc.com with title [企業溝通無國界翻譯系統] C* after meeting.**

# Submit Results

## Transcript & Proper Nouns

- **Submit the transcript & proper nouns list you identified based on the testing wav, upload to <span style="color:red">your own bucket before 2/15 13:00</span>**

- **List all detected proper nouns. If there are duplicates, don't list only one.**

- **Naming by [Cnumber_transcript] and [Cnumber_proper_nouns].**
  - •C1_ transcript.txt
  - •C2_ proper_nouns.txt

大家好，今天要討論的是關於DDR Ratio的問題，在dp上發現這週的ratio很高，請問MARTIN是否知道發生原因 ？
很抱歉我昨天值大夜班，有把事情交接給LISA了，可以請他說明原因。
關於這周DDR ratio過高的原因可能是EC被動過的原因，我回去和母版比對後發現溫度等數值都不太一樣。
為什麼EC會被更改過，數值不是應該和母版對齊嗎? IT能不能查一下系統的log，確認做change的人是誰?
可以，我回去撈一下資料。
另外，IT能否也將做change的資料上架到dp，當有人做了不符合權限的事情可以印出資料，並自動寄送alert信件給相關人員。
好的，這件事技術上沒問題，但我需要回去和我老闆討論一下，因為這屬於架構上的change，我這邊需要新增cloud function來抓log的資料，BigQuery那邊也需要新增table欄位才行。
好，那請你下次再update這件事給我。另外EC被動過這件事也請Martin追一下發生原因，也請你下次update給我，謝謝。
好的，我這邊會持續追蹤這件事。
好了今天的會議就開到這邊，謝謝大家。
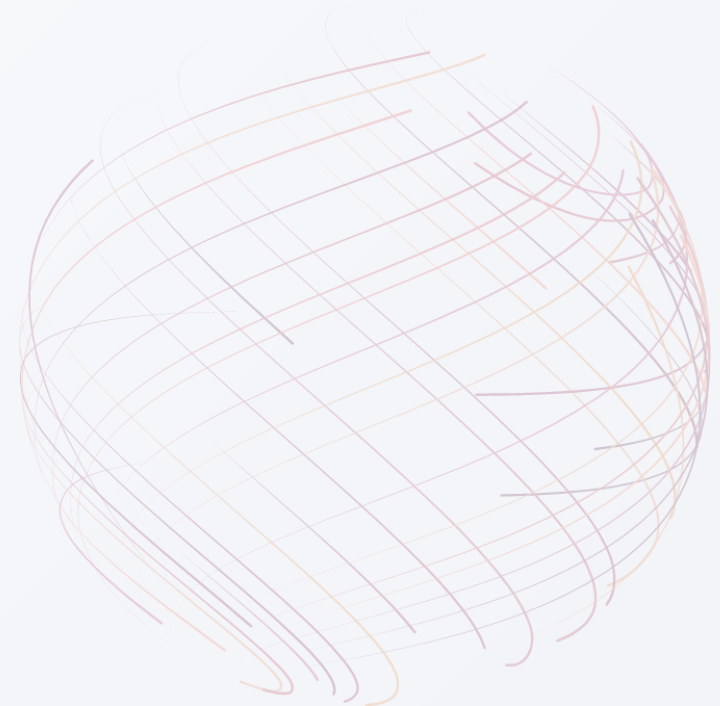謝謝。
謝謝。
掰掰。

C1_ transcript.txt

DDR Ratio
dp
大夜
DDR Ratio
EC
EC
dp
cloud function
BigQuery
EC

C2_ proper_nouns.txt

2025 | TSMC IT **CareerHack**
Disruptive Innovation in AI Era

# Submit Results

## Demo ppt & Presentation

- **Upload your demo ppt to your own bucket before 2/15 13:45.**

- **Demo ppt naming by [Cnumber_vnumber].**
  - C1_v1.pptx
  - C2_v3.pptx

- **Presentation Timeline Description:**
  - Presentation: 15 mins (include live demo)
  - QA: 5 mins
  - Live demo: 1 min

# Timeline

**2/7 Workshop**

1. Workshop
2. Provide cloud resources
3. Provide training dataset

**2/14 Hack Day 1**

1. Hack Day 1
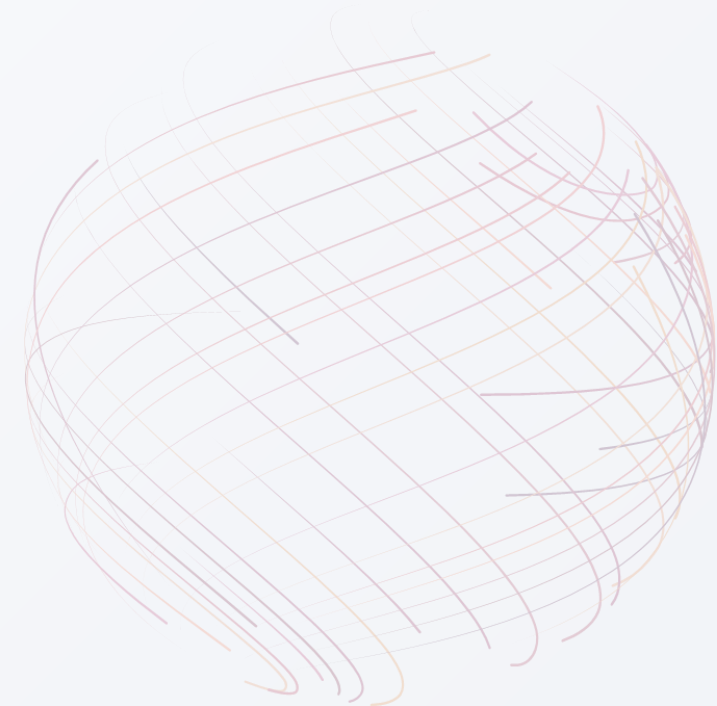2. Provide testing dataset

**2/15 Hack Day 2**

1. Hack Day 2
2. Submit the transcript and the proper nouns you identified based on the testing wav, upload to your own bucket before 13:00
3. Present your project
4. Upload your ppt before 13:45

# Reference

**GCP Vertex AI RAG Sample Code: [Link](#)**

**GCP Speech-to-text Sample Code: [Link](#)**

**Langchain with Cloud SQL: [Link](#)**

Q&A

2025 TSMC IT
CareerHack

# Thank You

SEE YOU NEXT PRESENTATION