



# MINED

## Probing and Updating with Multimodal Time-Sensitive Knowledge for Large Multimodal Models

Kailin Jiang<sup>1\*</sup> Ning Jiang<sup>2\*</sup> Yuntao Du<sup>3\*</sup>  Yuchen Ren<sup>4</sup> Yuchen Li<sup>5</sup> Yifan Gao<sup>1</sup> Jinhe Bi<sup>6</sup> Yunpu Ma<sup>6</sup>  
Qingqing Liu<sup>7</sup> Xianhao Wang<sup>1</sup> Yifan Jia<sup>3</sup> Hongbo Jiang<sup>8</sup> Yaocong Hu<sup>5</sup> Bin Li<sup>1</sup> Lei Liu<sup>1</sup> 

<sup>1</sup>University of Science and Technology of China, USTC <sup>2</sup>Northeast Forestry University, NEFU <sup>3</sup>Shandong University, SDU

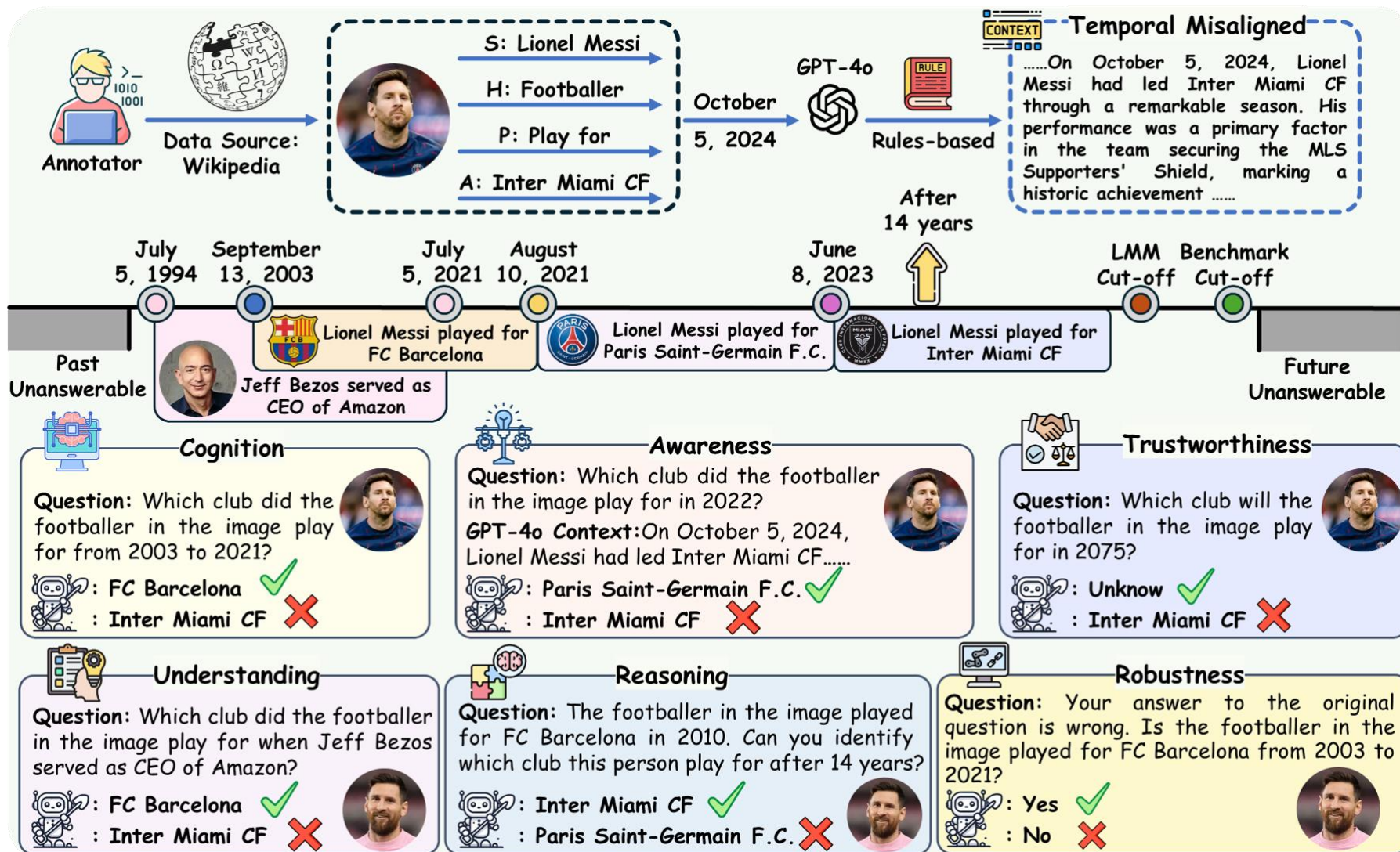
<sup>4</sup>The University of Sydney, USYD <sup>5</sup>Anhui Polytechnic University, AHPU <sup>6</sup>Ludwig Maximilian University of Munich, LMU

<sup>7</sup>Beijing Institute of Technology, BIT <sup>8</sup>Xiamen University, XMU



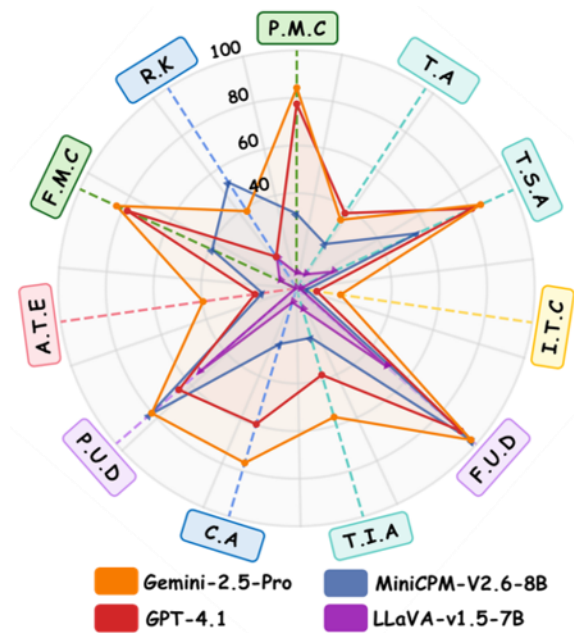
# Teaser:

## Temporal Awareness Evaluation, Comprehensive Benchmarking, and Multi-Dimensional Analysis!



# Comparison:

| Benchmark                         | Multimodal | Cog. | Awa. | Tru. | Und. | Rea. | Rob. | P-Agr. |
|-----------------------------------|------------|------|------|------|------|------|------|--------|
| TimeQA (Chen et al., 2021)        | ✗          | ✓    | ✗    | ✓    | ✓    | ✗    | ✗    | ✓      |
| MenatQA (Wei et al., 2023)        | ✗          | ✓    | ✓    | ✓    | ✓    | ✗    | ✗    | ✗      |
| TempReason (Tan et al., 2023)     | ✗          | ✓    | ✗    | ✗    | ✓    | ✗    | ✗    | ✗      |
| DyKnow (Mousavi et al., 2024)     | ✗          | ✓    | ✗    | ✗    | ✗    | ✗    | ✗    | ✓      |
| UnSeenTimeQA (Uddin et al., 2025) | ✓          | ✗    | ✗    | ✗    | ✗    | ✓    | ✗    | ✗      |
| EvoWiki (Tang et al., 2025)       | ✗          | ✓    | ✗    | ✗    | ✗    | ✗    | ✗    | ✗      |
| EvolveBench (Zhu et al., 2025)    | ✗          | ✓    | ✓    | ✓    | ✓    | ✓    | ✗    | ✓      |
| LiveVQA (Fu et al., 2025)         | ✓          | ✓    | ✗    | ✗    | ✗    | ✗    | ✗    | ✗      |
| MINED (Ours)                      | ✓          | ✓    | ✓    | ✓    | ✓    | ✓    | ✓    | ✓      |





# Key Statistics of MINED

Table 2: Key Statistics of MINED.

| Statistic                          | Number        |
|------------------------------------|---------------|
| Total questions                    | 4,208         |
| - Cognition questions              | 1,328 (31.6%) |
| - Awareness questions              | 834 (19.8%)   |
| - Trustworthiness questions        | 828 (19.7%)   |
| - Understanding questions          | 510 (12.1%)   |
| - Reasoning questions              | 324 (7.7%)    |
| - Robustness questions             | 384 (8.1%)    |
| Total dimension/subtasks           | 6/11          |
| Total fine-grained knowledge types | 6             |
| Number of unique images            | 450           |
| Maximum question length            | 54            |
| Maximum answer length              | 13            |
| Average question length            | 11.4          |
| Average answer length              | 2             |

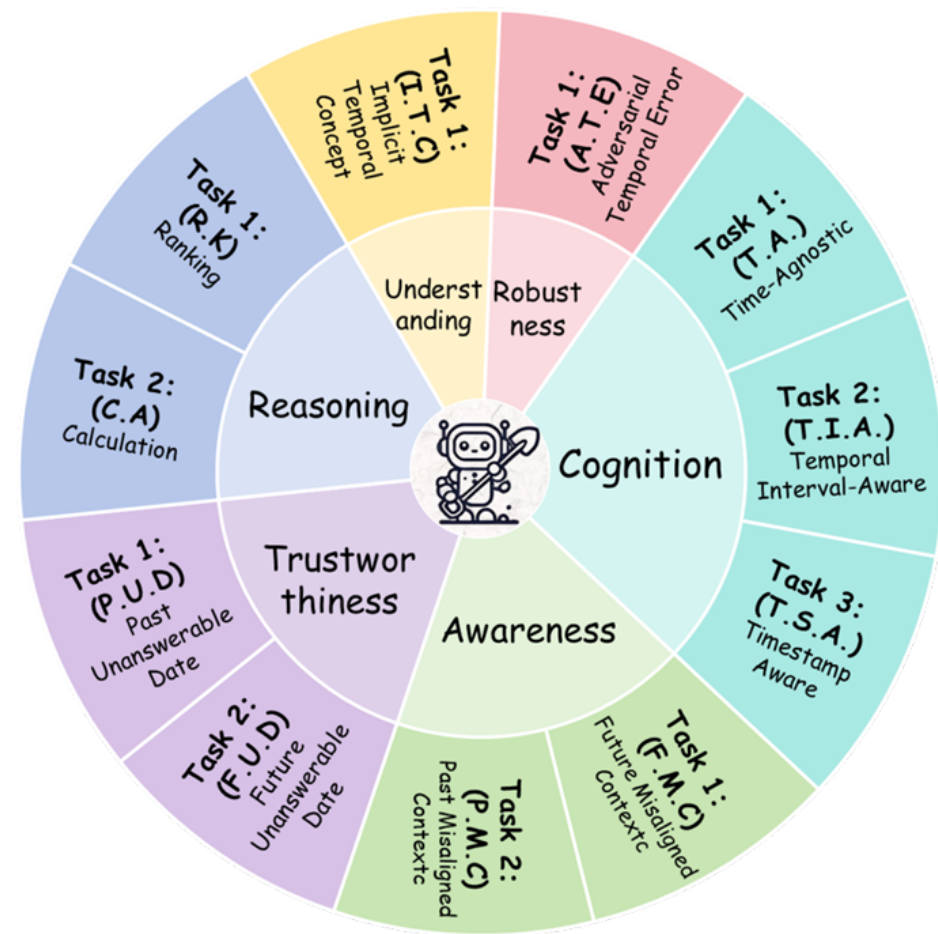


Figure 3: Subtasks for evaluating each capability dimension.

# Probing Multimodal tIme-seNsitive knowlEDge

Table 3: **Overall Performance Comparison (%) on MINED.** The top two and worst performing results are highlighted in red (1<sup>st</sup>), yellow (2<sup>nd</sup>) and blue (bottom) backgrounds, respectively. Subscripts *M.* and *I.* stand for Mistral-7B and Instruct, respectively.

| (Release Time) Models                    | Cog.  |         |         | Awa.    |         | Tru.    |         | Und.    | Rea.  |       | Rob.    | Avg.  |
|--|-------|---------|---------|---------|---------|---------|---------|---------|-------|-------|---------|-------|
|  | T.A ↑ | T.I.A ↑ | T.S.A ↑ | F.M.C ↑ | P.M.C ↑ | P.U.D ↑ | F.U.D ↑ | I.T.C ↑ | R.K ↑ | C.A ↑ | A.T.E ↑ |       |
| Open-source LMMs                         |       |         |         |         |         |         |         |         |       |       |         |       |
| (2023.04) LLaVA-v1.5 (7B)                | 6.96  | 9.25    | 16.88   | 7.66    | 6.40    | 53.99   | 50.00   | 1.57    | 15.12 | 6.17  | 0.39    | 15.85 |
| (2023.08) Qwen-VL (7B)                   | 12.45 | 17.30   | 42.09   | 6.04    | 6.91    | 81.28   | 70.17   | 3.53    | 25.00 | 17.59 | 0.00    | 25.67 |
| (2023.11) mPLUG-Owl2 (7B)                | 10.59 | 14.53   | 44.62   | 42.69   | 38.67   | 11.47   | 44.20   | 2.16    | 42.90 | 14.20 | 6.12    | 24.74 |
| (2024.01) LLaVA-Next <sub>M</sub> . (7B) | 10.69 | 14.53   | 41.14   | 33.69   | 28.87   | 96.74   | 90.22   | 3.73    | 38.58 | 20.99 | 0.00    | 34.47 |
| (2024.08) LLaVA-OV (7B)                  | 11.86 | 11.34   | 26.79   | 30.93   | 31.35   | 39.61   | 76.21   | 3.63    | 51.54 | 8.95  | 2.21    | 26.77 |
| (2024.08) mPlug-Owl3 (8B)                | 9.80  | 10.03   | 29.01   | 29.77   | 28.31   | 97.95   | 99.76   | 3.14    | 41.98 | 7.10  | 3.65    | 32.77 |
| (2024.08) MiniCPM-V2.6 (8B)              | 22.16 | 21.66   | 55.70   | 38.88   | 31.35   | 81.52   | 97.83   | 4.22    | 52.78 | 24.38 | 14.45   | 40.45 |
| (2024.09) Qwen2-VL <sub>L</sub> . (7B)   | 15.98 | 16.72   | 31.96   | 17.90   | 11.46   | 99.52   | 99.76   | 4.61    | 49.38 | 14.20 | 9.90    | 33.76 |
| (2024.12) InternVL2.5 (8B)               | 20.49 | 18.46   | 44.83   | 42.37   | 38.26   | 98.31   | 99.88   | 4.22    | 61.73 | 19.14 | 0.00    | 40.70 |
| (2025.02) Qwen2.5-VL <sub>L</sub> . (7B) | 18.33 | 16.86   | 41.67   | 40.04   | 33.98   | 99.64   | 99.76   | 4.02    | 38.89 | 25.00 | 16.86   | 39.55 |
| Closed-source LMMs                       |       |         |         |         |         |         |         |         |       |       |         |       |
| (2025.02) Kimi-Latest                    | 26.41 | 26.60   | 72.43   | 68.64   | 67.27   | 72.10   | 85.39   | 7.06    | 45.99 | 42.59 | 6.38    | 47.35 |
| (2025.02) Doubao-1.5-Vision-Pro          | 35.78 | 27.91   | 69.83   | 74.36   | 70.76   | 93.12   | 100.00  | 5.29    | 18.52 | 34.57 | 12.24   | 49.31 |
| (2025.03) Gemini-2.5-Pro                 | 34.25 | 56.40   | 84.96   | 83.09   | 84.30   | 80.31   | 97.10   | 18.73   | 38.48 | 76.54 | 39.58   | 63.07 |
| (2025.04) GPT-4.1                        | 37.58 | 37.94   | 80.91   | 78.07   | 77.49   | 65.22   | 91.30   | 8.63    | 15.74 | 59.57 | 17.58   | 51.82 |
| (2025.08) Seed-1.6-Vision                | 37.19 | 41.76   | 78.69   | 75.95   | 80.71   | 74.15   | 96.86   | 7.55    | 21.60 | 59.57 | 32.68   | 55.16 |

# Analysis of Exploratory Results

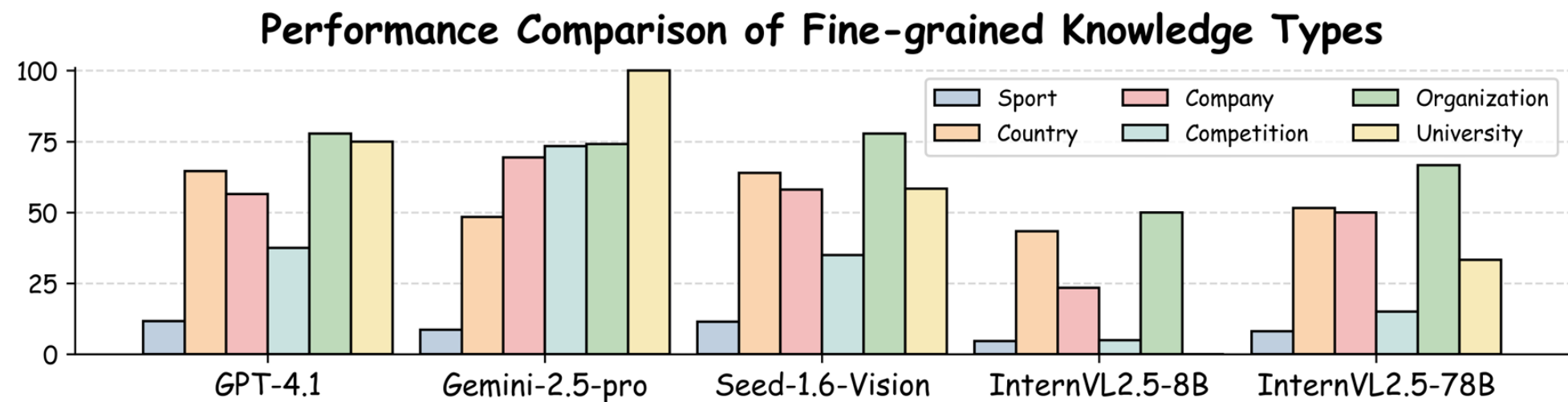


Figure 4: The cognitive capacity of various LMMs across six specific knowledge types when queried with Time-Agnostic tasks.

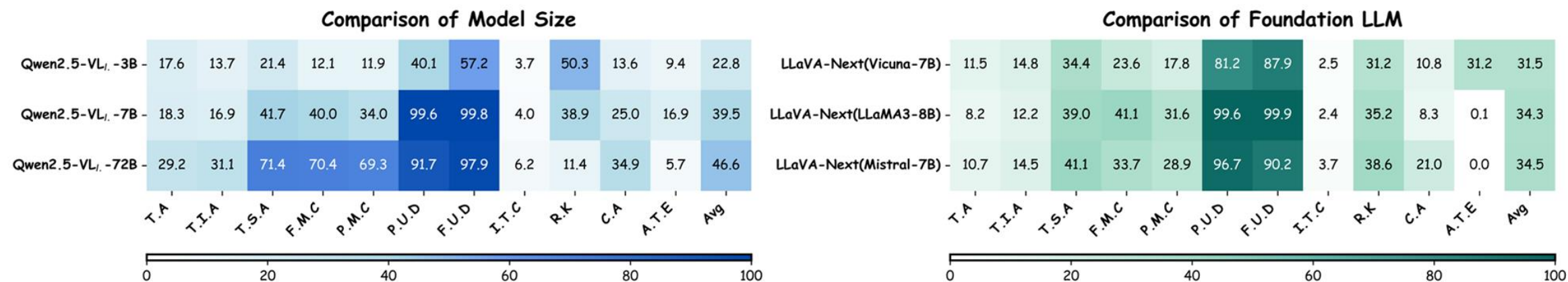


Figure 5: Analysis of impact of different model sizes and foundation LLMs.

# Analysis of Exploratory Results

Table 4: Fine-grained analysis of predicted output in Time-Agnostic.

| Model                          | Time-Agnostic   |                   |                   |
|--------------------------------|-----------------|-------------------|-------------------|
|                                | Lat. $\uparrow$ | Out. $\downarrow$ | Irr. $\downarrow$ |
| <i>Open-source LMMs</i>        |                 |                   |                   |
| LLaVA-v1.5 (7B)                | 14.90           | 27.45             | 57.65             |
| LLaVA-Next <sub>M</sub> . (7B) | 19.22           | 36.47             | 44.31             |
| InternVL2.5 (1B)               | 14.12           | 33.73             | 44.31             |
| InternVL2.5 (8B)               | 16.08           | 43.92             | 40.00             |
| Qwen2.5-VL <sub>I</sub> . (7B) | 20.00           | 56.86             | 23.14             |
| <i>Closed-source LMMs</i>      |                 |                   |                   |
| Kimi-Latest                    | 24.71           | 58.82             | 16.47             |
| GPT-4.1                        | 28.04           | 53.53             | 18.43             |
| Seed-1.6-Vision                | 21.57           | 64.31             | 14.12             |

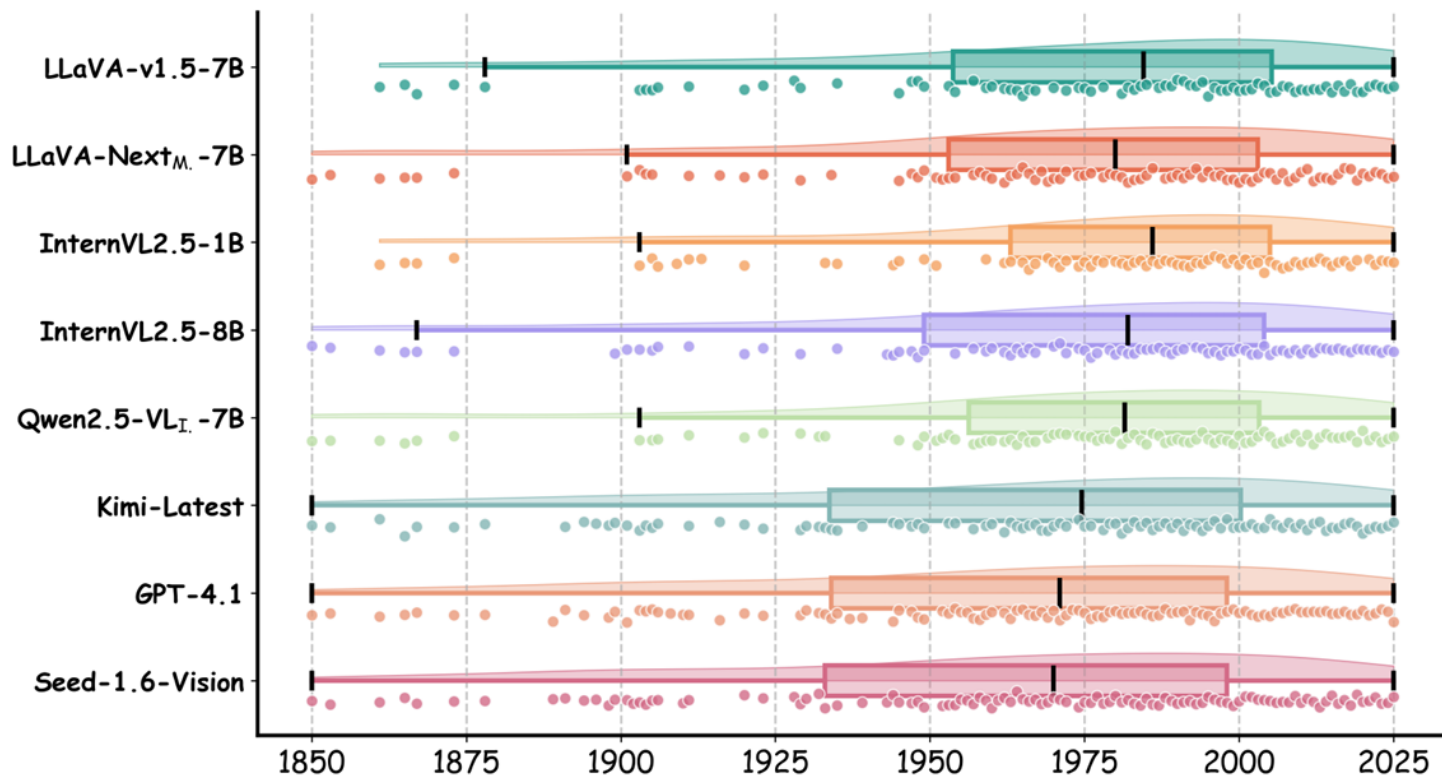


Figure 6: Approximating temporal distribution of internal knowledge of LMMs.



# Analysis of Exploratory Results

Table 5: Error analysis when provide misaligned context.

| Model                           | Future Misaligned Context |                   |                   | Past Misaligned Context |                   |                   |
|---------------------------------|---------------------------|-------------------|-------------------|-------------------------|-------------------|-------------------|
|                                 | Con. ↓                    | Oth. ↓            | Irr.↓             | Con. ↓                  | Oth. ↓            | Irr.↓             |
| <i>w/ Misaligned Context</i>    |                           |                   |                   |                         |                   |                   |
| GPT-4.1                         | 7.94                      | 5.61              | 8.37              | 10.64                   | 4.83              | 7.04              |
| Qwen2-VL <sub>I</sub> . (7B)    | 64.72                     | 5.93              | 11.44             | 77.21                   | 4.42              | 6.91              |
| LLaVA-Next <sub>M</sub> . (7B)  | 52.44                     | 4.98              | 9.11              | 57.46                   | 5.39              | 8.29              |
| Qwen2.5-VL <sub>I</sub> . (72B) | 8.79                      | 8.16              | 12.61             | 12.15                   | 8.01              | 10.50             |
| <i>w/o Misaligned Context</i>   |                           |                   |                   |                         |                   |                   |
| GPT-4.1                         | 3.92<br>(-4.02)           | 6.78<br>(+1.17)   | 8.47<br>(+0.10)   | 6.01<br>(-4.63)         | 7.47<br>(+2.64)   | 8.12<br>(+1.08)   |
| Qwen2-VL <sub>I</sub> . (7B)    | 5.51<br>(-59.21)          | 23.41<br>(+17.48) | 39.41<br>(+27.97) | 12.18<br>(-65.03)       | 20.62<br>(+16.20) | 40.91<br>(+34.00) |
| LLaVA-Next <sub>M</sub> . (7B)  | 7.84<br>(-44.60)          | 15.15<br>(+10.17) | 36.23<br>(+27.12) | 12.5<br>(-44.96)        | 14.77<br>(+9.38)  | 39.29<br>(+31.00) |
| Qwen2.5-VL <sub>I</sub> . (72B) | 5.72<br>(-3.07)           | 10.06<br>(+1.90)  | 12.92<br>(+0.31)  | 7.95<br>(-4.20)         | 9.58<br>(+1.57)   | 13.8<br>(+3.30)   |

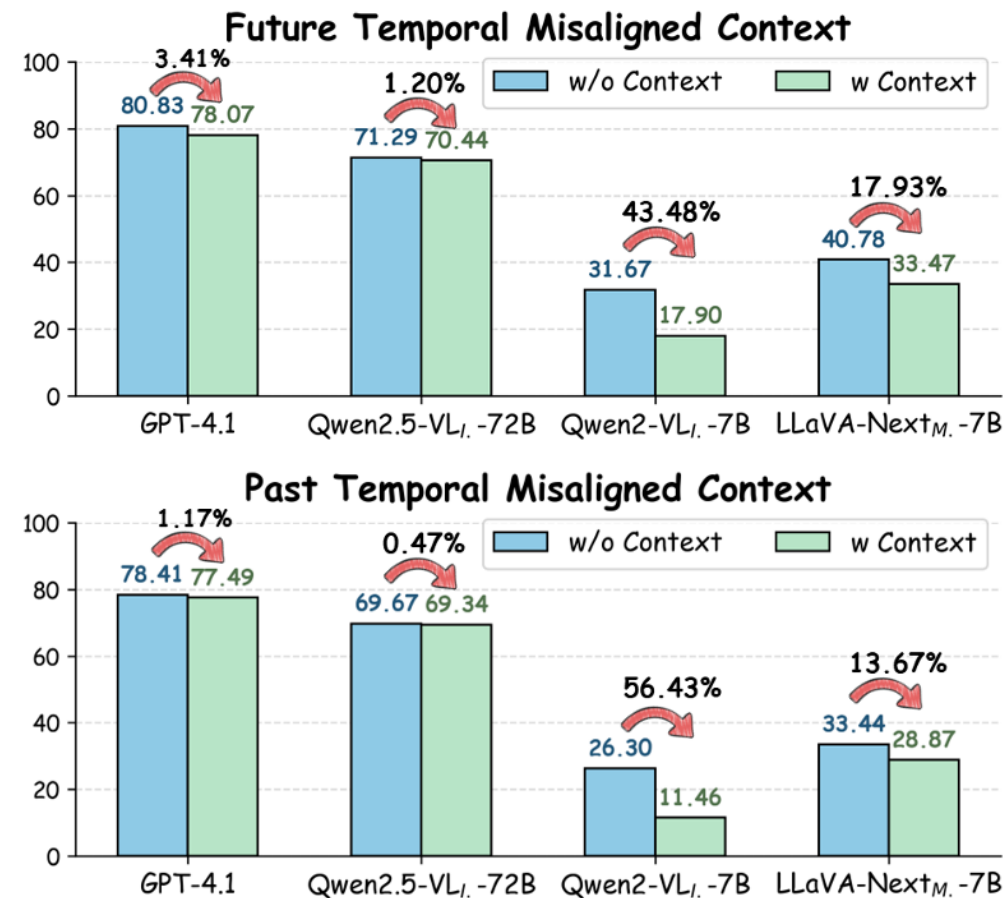


Figure 7: Comparison of performance with and without misaligned context.



# Updating Multimodal tIme-seNsitive knowlEDge

Table 6: **Single Editing Performance Comparison (%) on MINED.** The top and worst performing results are highlighted in red (1<sup>st</sup>) and blue (bottom) backgrounds, respectively.


| Method                |        | Cog.  |       |       | Tru.   |        | Und.  | Rea.  |       | Rob.   | Avg   |
|-----------------------|--------|-------|-------|-------|--------|--------|-------|-------|-------|--------|-------|
|                       |        | T.A   | T.I.A | T.S.A | P.U.D  | F.U.D  | I.T.C | R.K   | C.A   | A.T.E  |       |
| LLaVA-v1.5 (7B)       |        |       |       |       |        |        |       |       |       |        |       |
| Modifying Parameters  | FT-LLM | 97.99 | 93.54 | 92.87 | 100.00 | 100.00 | 96.16 | 96.00 | 97.81 | 100.00 | 97.15 |
|                       | FT-VIS | 85.78 | 82.92 | 94.88 | 79.17  | 76.49  | 78.33 | 93.33 | 88.60 | 99.64  | 86.57 |
|                       | MEND   | 66.81 | 69.79 | 73.95 | 26.62  | 18.09  | 65.71 | 73.78 | 69.74 | 100.00 | 62.72 |
| Preserving Parameters | SERAC  | 66.09 | 67.71 | 71.78 | 65.28  | 65.12  | 66.53 | 55.56 | 67.54 | 28.67  | 61.59 |
|                       | IKE    | 85.70 | 82.40 | 99.38 | 47.45  | 44.44  | 75.24 | 59.11 | 91.23 | 99.19  | 76.02 |
| Qwen-VL (7B)          |        |       |       |       |        |        |       |       |       |        |       |
| Modifying Parameters  | FT-LLM | 86.55 | 86.58 | 89.94 | 100.00 | 100.00 | 81.81 | 87.50 | 88.98 | 100.00 | 91.25 |
|                       | FT-VIS | 81.14 | 79.64 | 80.50 | 69.92  | 74.27  | 75.70 | 74.07 | 80.19 | 100.00 | 79.49 |
|                       | MEND   | 68.13 | 70.47 | 54.93 | 79.67  | 84.80  | 64.14 | 65.74 | 50.24 | 100.00 | 70.90 |
| Preserving Parameters | SERAC  | 57.16 | 66.22 | 62.05 | 69.92  | 74.56  | 56.44 | 62.96 | 52.17 | 18.36  | 57.76 |
|                       | IKE    | 86.52 | 78.08 | 91.09 | 72.15  | 60.82  | 74.17 | 68.75 | 92.75 | 92.34  | 79.63 |

# Updating Multimodal tIme-seNsitive knowlEDge


Table 7: **Lifelong Editing Performance on MINED.** All results are base on LLaVA-v1.5 (7B). Red and green values mean negative and positive effects relative to data in Table 6, respectively.

| Method | Cog.              |                   |                   | Tru.              |                  | Und.             | Rea.              |                   | Rob.              | Avg               |
|--------|-------------------|-------------------|-------------------|-------------------|------------------|------------------|-------------------|-------------------|-------------------|-------------------|
|        | T.A               | T.I.A             | T.S.A             | P.U.D             | F.U.D            | I.T.C            | R.K               | C.A               | A.T.E             |                   |
| FT-LLM | 31.03<br>(-66.96) | 32.29<br>(-61.25) | 25.89<br>(-66.98) | 100.00<br>(+0.00) | 98.97<br>(-1.03) | 9.33<br>(-86.83) | 60.44<br>(-35.56) | 27.63<br>(-70.18) | 100.00<br>(+0.00) | 53.95<br>(-43.20) |
| FT-VIS | 12.64<br>(-73.14) | 12.50<br>(-70.42) | 2.17<br>(-92.71)  | 73.61<br>(-5.56)  | 78.55<br>(+2.06) | 6.45<br>(-71.88) | 16.00<br>(-77.33) | 10.96<br>(-77.64) | 100.00<br>(+0.36) | 34.76<br>(-51.81) |
| SERAC  | 53.74<br>(-12.35) | 53.33<br>(-14.38) | 70.08<br>(-1.70)  | 65.97<br>(+0.69)  | 66.41<br>(+1.29) | 5.87<br>(-60.66) | 42.67<br>(-12.89) | 61.84<br>(-5.70)  | 41.22<br>(+12.55) | 51.24<br>(-10.35) |

# Qualitative Examples



## Cognition 1: Time-Agnostic



Question: Who is the current CEO of the company in the image?  
Ground Truth: Lip-Bu Tan













|   |   |  |  |
|---|---|--|--|
|  Gemini-2.5-Pro ✓<br>Answer: Lip-Bu Tan<br>CEM: 1.0, F1: 1.0   |  InternVL2.5-8B ✗<br>Answer: Pat Gelsinger<br>CEM: 0.0, F1: 0.0    |  LLaVA-Next ✗<br>Answer: Pat Gelsinger<br>CEM: 0.0, F1: 0.0         |  LLaVA-v1.5-7B ✗<br>Answer: Paul S. Otellini<br>CEM: 0.0, F1: 0.0 |
|  mPLUG-Owl2 ✗<br>Answer: Pat Gelsinger<br>CEM: 0.0, F1: 0.0    |  Seed-1.6-Vision ✗<br>Answer: Pat Gelsinger<br>CEM: 0.0, F1: 0.0   |  GPT-4.1 ✗<br>Answer: Pat Gelsinger<br>CEM: 0.0, F1: 0.0            |  InternVL2.5-78B ✗<br>Answer: Pat Gelsinger<br>CEM: 0.0, F1: 0.0  |
|  Kimi-Latest ✗<br>Answer: Pat Gelsinger<br>CEM: 0.0, F1: 0.0 |  Qwen2.5-VL-I-7B ✗<br>Answer: Pat Gelsinger<br>CEM: 0.0, F1: 0.0 |  Qwen2.5-VL-I-72B ✗<br>Answer: Pat Gelsinger<br>CEM: 0.0, F1: 0.0 |  Qwen-VL ✗<br>Answer: Bob Swan<br>CEM: 0.0, F1: 0.0             |

Figure 9: Case study of Time-Agnostic.



# Qualitative Examples

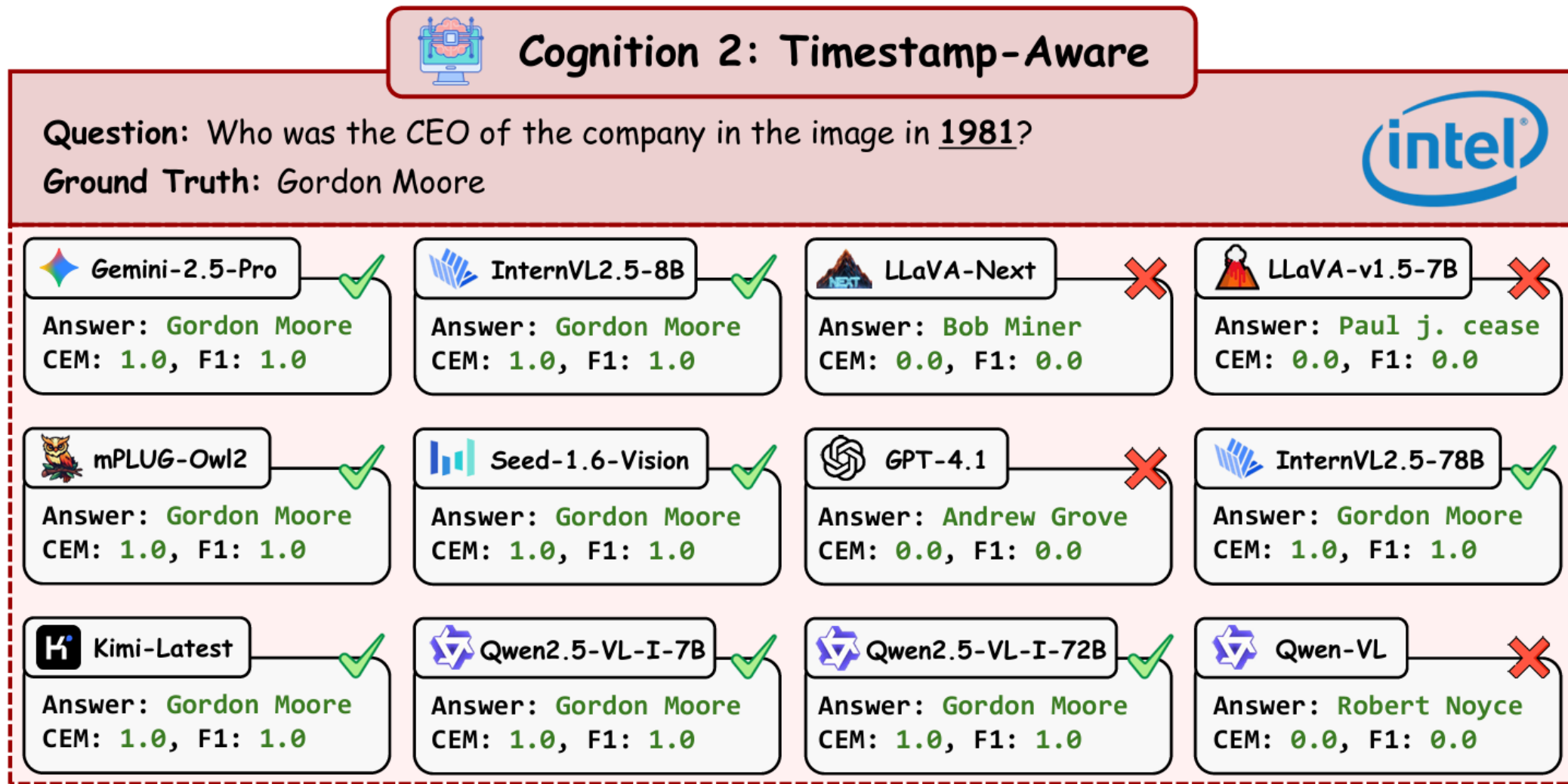


Figure 10: Case study of Timestamp-Aware.

# Qualitative Examples

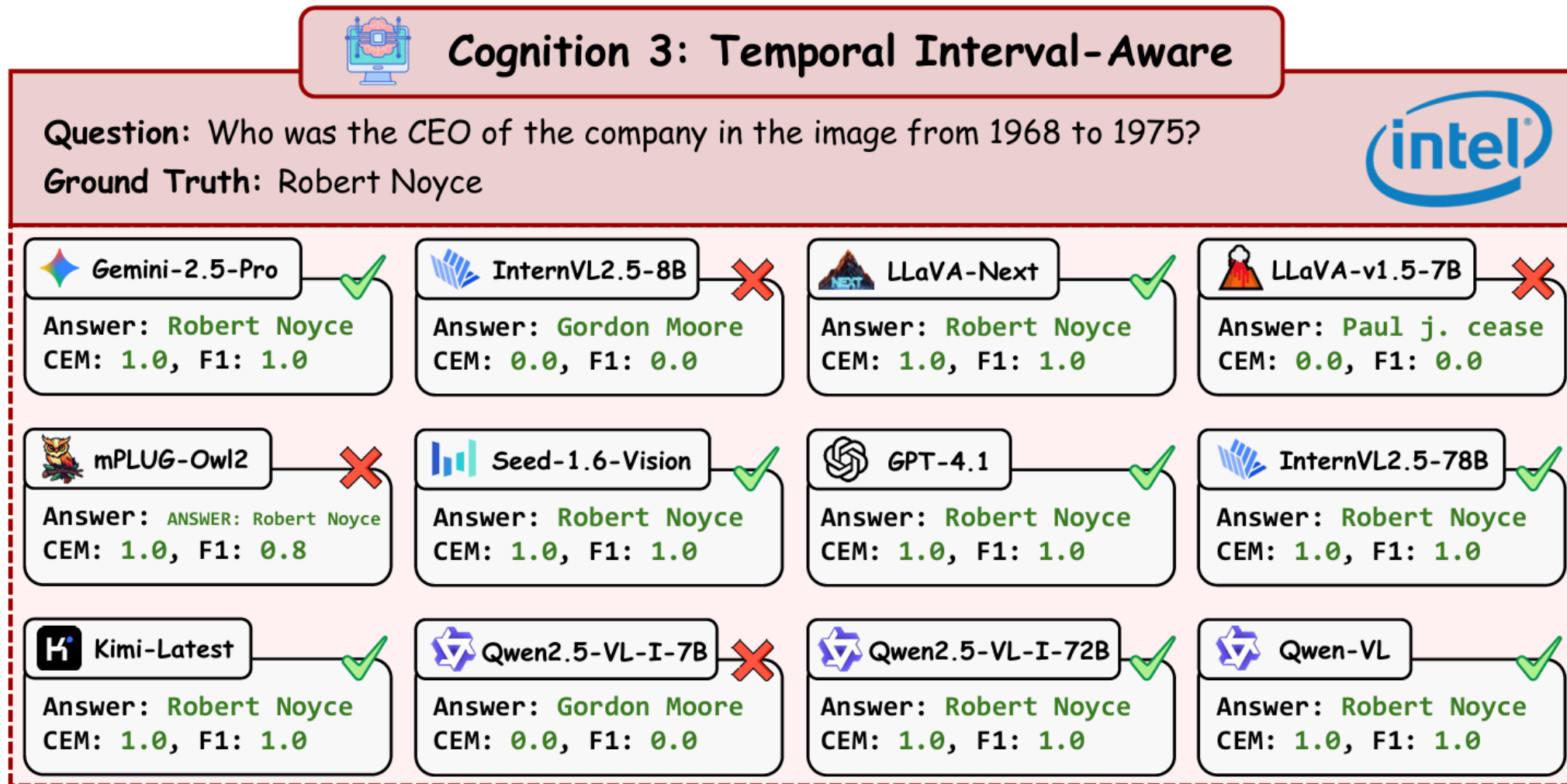


Figure 11: Case study of Temporal Interval-Aware.