# SynTab-LLaVA: Enhancing Multimodal Table Understanding with Decoupled Synthesis

## Supplementary Material

## A. Examples of MTU Datasets

We visualize the samples of the existing MTU datasets mentioned in Sec. 2.1, with the results shown in Fig. 1 and Fig. 2. The first to fourth columns in the figures represent the dataset name, the input table image, the user-provided question, and the expected answer, respectively.

For the Table Question Answering (TQA), the datasets WTQ [30], FeTaQA [27], HiTab [5], AIT-QA [17], TabMCQ [16], TAT-QA [51], and TABMWP [24] are visualized in Fig. 1. We observe that the requirements for answers vary across datasets. For example, WTQ only requires a concise response, while FeTaQA demands a coherent and lengthy textual answer generated based on the table content and the given question. HiTab and AIT-QA have a large number of table cells and complex table structures, yet they also only require a short answer. For TABMWP and TAT-QA, which focus on table-based numerical reasoning, the task involves generating a series of calculation steps followed by the final answer. For TabMCQ, the model only needs to match the table content with the options provided in the question and select the correct answer.

For Table Fact Verification (TFV), there are three datasets: TabFact [3], InfoTabs [13], and PubhealthTab [1]. This task is relatively straightforward, requiring the model to index, query, and reason over the table image based on the users input statement. If the derived conclusion aligns with the users statement, the output is "correct"; if the conclusion contradicts the users statement, the output is "incorrect"; and if it is impossible to determine the validity of the users statement solely based on the table content, the output is "neutral". The visualization of specific samples is shown in Fig. 2.

The Table-to-Text Generation (T2T) task includes four datasets: ToTTo [29], HiTab_T2T [5], Rotowire [38], and WikiBIO [18]. Examples from each dataset can be found in Fig. 2. These datasets require the model to generate coherent sentences as responses based on the users input question and the table image. The generated responses may summarize highlighted cells in the table image or provide a comprehensive description of a related biography.

The Table Structure Understanding (TSU) [46] task includes six sub-tasks, which are designed to evaluate the model's capability to extract basic structural information and understand table content. Table Size Detection (TSD) predicts the number of rows and columns in a given table. Table Cell Extraction (TCE) involves extracting the content of a target cell based on the row and column indices spec-

ified by the user. Table Cell Locating (TCL), on the other hand, requires identifying the positions of cells whose content matches the user-provided text. Merge Cell Detection (MCD) aims to detect merged cells and output their coordinates, including the top-left and bottom-right corners. Row & Column Extraction (RCE) focuses on extracting the content of entire rows or columns. Finally, Table Recognition (TR) requires the model to convert the input table image into corresponding HTML, LaTeX, or Markdown.

## B. Detailed Category Definition

In this part, we provide a detailed explanation of the 6 main question categories and 11 subcategories defined earlier.

### B.1. Retrieval

**Data retrieval** focuses on extracting specific information or multiple data points from tables, enhancing the model's ability to locate and return exact values based on a query. For instance, in MTU, data retrieval helps identify and extract relevant cells or ranges, directly improving the models performance in retrieving precise information from structured data.

### B.2. Data Operations

**Data counting** involves determining the number of occurrences of specific values or conditions within a table. This requires scanning through rows and columns to count how many times a particular criterion is met.

**Data order** includes tasks related to organizing or identifying specific statistical properties of data within a table. This involves: maximum, minimum, sorting, median, comparison of values, and so on. These tasks mainly focus on how to systematically order and interpret data based on its statistical attributes.

**Determine range** involves calculating the difference between two data points within a specific column or data set in a table. This could be the difference between the maximum and minimum values, or it could be between the third highest and minimum values, depending on the context.

**Filter operations** include questions related to selecting and manipulating subsets of data based on specific criteria. The key operations involve: filtering, union, intersection, set difference, subset, conditional filtering, *etc*.

### B.3. Numerical Calculations

**Numerical calculations** involve: addition: summing numerical values; subtraction: finding the difference between

| Dataset | Image | Question | Answer |
| --- | --- | --- | --- |
| WTQ |  | This image shows a table about 'Taekwondo at the 2010 South American Games', answer the following question using a single word or phrase. The final result should be presented in the JSON format of {"answer": [<a list of answer strings>]}. Question: who won the most gold medals? Answer: | Using the table as reference, the answer is {"answer": ['Brazil']}. |
| FeTaQA |  | Give you a screenshot of a table titled 'Personal bests' from the Wikipedia of 'Lauren Fleshman'. Provide a well-considered answer to the following question: When was Lauren Fleshman's personal best in the marathon achieved? | In 2011, Lauren Fleshman achieved her personal best in the marathon in New York City. |
| HiTab |  | Based on the excel table, answer the following question in brief. Your answer should be in the JSON structure, formatted as {"answer": [<a list of answer strings>]}. For instance, {"answer": ["2021-12-04", "72.33"]}.Table title: career statistics. Question: how many league goals did Hristo Yanev score for cska sofia? | Based on the table information, the final answer is {"answer": ["54.0"]}. |
| AIT-QA |  | What was United's cost per ASM in 2017? Give me a concise answer based on the given table about an airline company. Show your final answer in the JSON format {"answer": "<YOUR ANSWER>"}. | Using the table as reference, the answer is {"answer": "15.72"}. |
| TabMCQ |  | Select the correct answer option according to the table and the question. Provide the final answer in the JSON structure, using the format {"answer": "<YOUR ANSWER>"}. Question: Which of the following is artificial? Options: 1. Air 2. Water 3. Humans 4. House | According to the given table, the answer is {"answer": "4. House"}. |
| TABMWP |  | Based on the table, solve the math word problem. Give a solution step by step. Provide a concluding answer in a JSON structure, using the format {"answer": "<YOUR ANSWER>"}. Problem: Victor bought 3 kilograms of amethyst and 2 kilograms of fool's gold. How much did he spend? (Unit: $) Solution: | Find the cost of the amethyst: $7.61 × 3 = $22.83. Find the cost of the fool's gold: $8.54 × 2 = $17.08. Now the total cost: $22.83 + $17.08 = $39.91. Therefore, the final answer is {"answer": "39.91"}. |
| TAT-QA |  | Read the table to answer the following questions in brief. Present the final answer in a JSON format {"answer": [<a list of answer strings>]}. How much is the 2019 salaries and fees? | Based on the table information, the answer is {"answer": ["4"]}. |
| RCE |  | Provide the contents of the cells located in row 4 of this table. Show the cells in a row in the JSON format {"row_id":"<row index>", "cell_list": "<a list of cells in this row>"}. | The target cell values of row 4: {"row_id": "4", "cell_list": ['Publisher', 'Nelson Thornes, 2000']} |
| TR |  | This picture illustrates a table. Please represent this table with the markdown-format in text. | \| \|For the Years Ended December 31, \| \| \|\n\| --- \| --- \| --- \| --- \|\n\| 2019 \| 2018 \| Increase (Decrease) \|\n\| Net bookings \| $6,388 \| $7,262 \| $(874) \|\n\| In-game net bookings \| $3,366 \| $4,203 \| $(837) \| |

Figure 1. Examples of existing MTU datasets samples.

# Dataset / Image / Question / Answer

## TabFact

| rank | nation | gold | silver | bronze | total |
|---|---|---|---|---|---|
| 1 | russia | 4 | 1 | 1 | 6 |
| 2 | united states | 0 | 1 | 1 | 2 |
| 3 | china | 0 | 1 | 0 | 1 |
| 3 | france | 0 | 1 | 0 | 1 |
| 5 | canada | 0 | 0 | 1 | 1 |
| 5 | poland | 0 | 0 | 1 | 1 |

**Question:** In this task, the goal is to distinguish whether the following sentence is confirmed or refuted by the table. Format your final answer as a JSON, using the structure {"answer\": "<YOUR ANSWER>"}.
poland took home the most gold medals

**Answer:** Utilizing the table as evidence, we can find that the sentence is false. So, the answer is {"answer": "refuted"}.

## InfoTabs

**Crufts**

| Formation | 1891 |
|---|---|
| Type | Dog show |
| Headquarters | Birmingham, England |
| Location | • United Kingdom |
| Official language | English |
| Website | www.crufts.org.uk |

**Question:** Refer to the table to decide if the following claim is affirmed, disproved by the given table, or if unclear, regard it as 'not enough information'. Format your final answer as a JSON, using the structure {"answer": "<YOUR ANSWER>"}. The Crufts dog show has been around under twenty years.

**Answer:** Based on the evidence in the table, the claim is not true. Thus, the final answer is {"answer": "disproved"}.

## PubhealthTab

| Party | As of June 2021 |
|---|---|
| Democratic Party | 41 |
| Republican Party | 18 |
| Vacancies | 0 |
| Total | 59 |

**Question:** claim: More women are restrict of domestic violence. Review the table named 'Trafficking Routes' to determine if the claim is supported or opposed by the table. Use 'not enough information' for inconclusive evidence. Show your final answer in the JSON format {"answer": "<YOUR ANSWER>"}

**Answer:** The information in the table is not sufficient to judge whether the claim is true or false. Therefore, the final answer: {"answer": "not enough information"}.

## ToTTo

| | Amine Scrubbing | Ca-Looping |
|---|---|---|
| Cost/CO 2 avoided | ~ $35–96/ton | ~$23.70/ ton |
| Raw material cost | $1,250/ton MEA | $25/ton CaCO3 |
| Efficiency Penalty | 6-12% | 6-8% |

**Question:** Presented is a Wikipedia table named 'Calcium looping', which is extracted from a section heading 'Cost per metric ton for CO2 captured'. As you can see, some table cells are highlighted with a different color. Please generate a one-sentence summary of these cells.

**Answer:** The cost of CO2 capture is $35 to $96 per ton for Amine Scrubbing and $23 to $70 per ton for Ca-Looping.

## HiTab_T2T

| sub-groups of the agri-food industry | eastern ontario | | northern ontario | |
|---|---|---|---|---|
| | french-language workers | other workers | french-language workers | other workers |
| | percent | | | |
| input and service supply | 2.0 | 2.1 | 2.0 | 1.3 |
| food, beverage, and tobacco processing | 8.7 | 6.0 | 3.0 | 3.3 |
| food retail and wholesale | 35.3 | 31.3 | 39.1 | 37.3 |
| food service | 52.1 | 60.6 | 55.0 | 58.1 |

**Question:** Please succinctly describe the highlighted cells within the given table titled 'agri-food industry sub-groups for workers aged 15 years and over, two agricultural regions of ontario, 2011'.

**Answer:** in eastern ontario, 60.6% of other agri-food workers worked in the restaurant and food services sector.

## WikiBIO

**Danny lloyd**

| name | Danny lloyd |
|---|---|
| birth_date | 1 january 1973 |
| birth_place | Chicago, illinois, u.s. |
| occupation | Teacher, actor |
| years_active | 1979 -- 1982 |

**Question:** Given the fact table in this image, which details information about 'Danny lloyd', create a succinct biography for this person.

**Answer:** Danny lloyd (born january 1, 1973) is an american teacher and former child actor.

## TSD

| | Tapping the ground | Raking the ground |
|---|---|---|
| 5-10 worms | 6 | 6 |
| 11-20 worms | 4 | 2 |

**Question:** I need to know the count of rows and columns in this specific table. Format your final answer as a JSON, using the structure {"row_number": "m", "column_number": "n"}.

**Answer:** This table has 3 rows and 3 columns. Thus, the final answer is {"row_number": "3", "column_number": "3"}.

## TCE

| age group of accused | females | |
|---|---|---|
| | number | rate |
| 12 to 17 | 12,273 | 1,096 |
| 18 to 24 | 16,845 | 1,043 |
| 25 to 34 | 23,481 | 922 |
| 35 to 44 | 14,478 | 596 |
| 45 to 54 | 8,846 | 353 |

**Question:** There is a cell located in the 5th row and the 1st column in this table, identify this cell and report its content. Provide the value of target cell in the JSON structure, using the format {"row_id":"m", "column_id": "n", "cell_value": "<Corresponding Cell Value>"}.

**Answer:** The target cell value in the 5th row and the 1st column is {"row_id": "5", "column_id": "1", "cell_value": "25 to 34"}.

## TCL

| round | player | position | nationality | college / junior / club team (league) |
|---|---|---|---|---|
| 1 | mark suzor | defense | canada | kingston canadians ( oha ) |
| 2 | drew callander | defense | canada | regina pats ( wchl ) |
| 3 | craig hanmer | defense | united states | mohawk valley comets ( nahl ) |
| 4 | dave hynek | defense | canada | kingston canadians (oha) |
| 5 | robin lang | defense | canada | cornell big red ( ecac ) |
| 6 | paul klasinski | left wing | united states | st paul vulcans (mjhl) |
| 7 | ray kurpis | right wing | united states | austin mavericks (mjhl) |

**Question:** Based on the table, what are the locations of the following cells?(1) 'craig hanmer'; Return the cell location in the JSON format {'value': '<cell value>', 'location': (row_id, column_id)}

**Answer:** The coordinates of target cells are listed below.
(1) {'value': 'craig hanmer', 'location': (4, 2)}

## MCD

| region of residence | number of syrian refugees | number of refugees from other countries |
|---|---|---|
| | | number |
| atlantic provinces | 2,670 | 765 |
| montreal census metropolitan area | 4,265 | 1,195 |
| rest of quebec | 1,030 | 1,970 |
| toronto census metropolitan area | 4,035 | 3,285 |
| rest of ontario | 6,175 | 2,910 |
| prairies | 4,735 | 4,055 |
| british columbia | 2,125 | 1,410 |
| canada | 25,035 | 15,580 |

**Question:** A merged cell in a table can be described using the JSON {'top-left':(A, B), 'bottom-right':(C, D)}. Here, '(A, B)' represents the top-left and '(C, D)' the bottom-right cells' (row_id, column_id) of the merged region. Using the table image, identify all such merged cells.

**Answer:** The merged cell regions are as follows: {'top-left': (1, 1), 'bottom-right':(2, 1)}; {'top-left':(2, 2), 'bottom-right':(2, 3)}.

Figure 2. Examples of existing MTU datasets samples.

## Our Rendered Images

| Warranty Reserve | Year Ended December 31, | | |
|---|---|---|---|
| | 2014 | 2013 | 2012 |
| Balance at the beginning of the year | $4,125 | $4,252 | $4,613 |
| Acquisition | 20,501 | 0 | 0 |
| Warranty expense | 12,909 | 7,440 | 6,828 |
| Warranty payments | (12,869) | (7,567) | (7,189) |
| Balance at the end of the period | $24,666 | $4,125 | $4,252 |

| (at December 31, 2007, in millions) | Amortized Cost | Fair Value |
|---|---|---|
| Due in one year or less | $3,288 | $3,292 |
| Due after 1 year through 5 years | 15,745 | 15,962 |
| Due after 5 years through 10 years | 17,558 | 17,777 |
| Due after 10 years | 20,453 | 20,773 |
| | 57,044 | 57,804 |
| Mortgage-backed securities | 7,108 | 7,116 |
| Total | $64,152 | $64,920 |

| | Year Ended December 31, | |
|---|---|---|
| | 2018 | 2017 |
| Reported GAAP Consolidated Net sales growth | 13.3% | 4.1% |
| Adjustments: | | |
| Impact of foreign currency translation(1) | (1.6)% | (0.6)% |
| Impact of Xplore acquisition(2) | (0.6)% | —% |
| Impact of Wireless LAN divestiture(3) | —% | 3.2% |
| Corporate, eliminations(4) | —% | (0.2)% |
| Consolidated Organic Net sales growth | 11.1% | 6.5% |

| Limited Partnership (Dollars in thousands) | Unfunded Commitments |
|---|---|
| SVB Strategic Investors Fund, LP | $6,793 |
| SVB Strategic Investors Fund II, LP | 25,143 |
| SVB Strategic Investors Fund III, LP | 128,093 |
| SVB Strategic Investors Fund IV, LP | 107,345 |
| SVB Capital Preferred Return Fund, LP | 41,256 |
| SVB Capital - NT Growth Partners, LP | 54,100 |
| Total | $362,730 |

| | Year Ended December 31, | | | |
|---|---|---|---|---|
| | | | $ | % |
| | 2008 | 2007 | Change | Change |
| | ($ in thousands) | | | |
| Distribution fees | | | | |
| U.S. high-grade | $30,287 | $34,939 | $(4,652) | (13.3)% |
| Eurobond | 14,143 | 8,148 | 5,995 | 73.6 |
| Total distribution fees | 44,430 | 43,087 | 1,343 | 3.1 |
| Variable transaction fees | | | | |
| U.S. high-grade | 16,260 | 17,602 | (1,342) | (7.6) |
| Eurobond | 4,003 | 10,680 | (6,677) | (62.5) |
| Other | 8,835 | 8,845 | (10) | (0.1) |
| Total transaction fees | 29,098 | 37,127 | (8,029) | (21.6) |
| Total commissions | $73,528 | $80,214 | $(6,686) | (8.3)% |

| | 2012 | | | |
|---|---|---|---|---|
| (In millions) | PrivateEquity | RealEstate | Other | Total |
| Beginning balance | $55 | $49 | $17 | $121 |
| Actual return on plan assets | 2 | - | - | 2 |
| Purchases | 12 | 10 | - | 22 |
| Sales | (13) | (5) | - | (18) |
| Ending balance | $56 | $54 | $17 | $127 |
| | 2011 | | | |
| (In millions) | PrivateEquity | RealEstate | Other | Total |
| Beginning balance | $46 | $37 | $17 | $100 |
| Actual return on plan assets | 7 | 5 | - | 12 |
| Purchases | 10 | 12 | - | 22 |
| Sales | (9) | (6) | - | (15) |
| Other | 1 | 1 | - | 2 |
| Ending balance | $55 | $49 | $17 | $121 |

| | Target Allocation | |
|---|---|---|
| | April 29, 2016 | April 24, 2015 |
| Asset Category | | |
| Equity securities | 49% | 49% |
| Debt securities | 23 | 23 |
| Other | 28 | 28 |
| Total | 100% | 100% |
| Non-U.S. Plans | | |
| | Target Allocation | |
| | April 29, 2016 | April 24, 2015 |
| Asset Category | | |
| Equity securities | 34% | 35% |
| Debt securities | 27 | 29 |
| Other | 39 | 36 |
| Total | 100% | 100% |

## Existing MTU Images

| Rank | Games | Player | Club | Career span |
|---|---|---|---|---|
| 1 | 432 | Brent Harvey | North Melbourne | 1996–2016 |
| 2 | 426 | Michael Tuck | Hawthorn | 1972–1991 |
| 3 | 403 | Kevin Bartlett | Richmond | 1965–1983 |
| 4 | 400 | Dustin Fletcher | Essendon | 1993–2015 |
| 5 | 383 | Robert Harvey | St Kilda | 1988–2008 |

| | 2019 | 2018 | 2017 |
|---|---|---|---|
| Restaurants sold to franchisees | — | 135 | 178 |
| New restaurants opened by franchisees | 19 | 11 | 18 |
| Proceeds from the sale of company-operated res... | | | |
| Cash (1) | $1,280 | $26,486 | $99,591 |
| Notes receivable | — | 70,461 | — |
| | $1,280 | $96,947 | $99,591 |
| Net assets sold (primarily property and equipm... | $— | $(21,329) | $(30,597) |
| Lease commitment charges (2) | — | — | (11,737) |
| Goodwill related to the sale of company-operat... | (2) | (4,663) | (10,062) |
| Other (3) | 88 | (24,791) | (9,161) |
| Gains on the sale of company-operated restaurants | $1,366 | $46,164 | $38,034 |

| Rank | Lane | Name | Nationality | Time | Notes |
|---|---|---|---|---|---|
| | 4 | Pieter van den Hoogenband | Netherlands | 48.30 | |
| | 3 | Alexander Popov | Russia | 48.69 | |
| | 7 | Gary Hall, Jr. | United States | 48.73 | |
| 4 | 5 | Michael Klim | Australia | 48.74 | |
| 5 | 2 | Neil Walker | United States | 49.09 | |
| 6 | 6 | Lars Frölander | Sweden | 49.22 | |
| 7 | 1 | Denis Pimankov | Russia | 49.36 | |
| 8 | 8 | Chris Fydler | Australia | 49.44 | |

| Party | Candidate | Votes | % | ± |
|---|---|---|---|---|
| Republican | Robert E. Smylie | 121,810 | 50.96% | |
| Democratic | Alfred M. Derr | 117,236 | 49.04% | |
| Majority | | 4,574 | | |
| Turnout | | 239,046 | | |
| Republican hold | | Swing | | |

| Particulars | Total | Male | Female |
|---|---|---|---|
| Total No. of Houses | 258 | | |
| Total Population | 1,396 | 715 | 681 |
| In the age group 0–6 years | 141 | 74 | 67 |
| Scheduled Castes (SC) | 984 | 510 | 474 |
| Scheduled Tribes (ST) | 0 | 0 | 0 |
| Literates | 912 | 498 | 414 |
| Illiterate | 484 | 217 | 267 |
| Total Worker | 471 | 399 | 72 |
| Main Worker | 305 | 274 | 31 |
| Marginal Worker | 166 | 125 | 41 |

Figure 3. Comparison of the table images synthesized in this work with existing MTU table images.

| | For the Years Ended December 31, | | | | | |
|---|---|---|---|---|---|---|
| Index | 2004 | 2005 | 2006 | 2007 | 2008 | 2009 |
| Aimco | 100.00 | 106.29 | 164.95 | 113.71 | 59.71 | 85.29 |
| MSCI US REIT | 100.00 | 112.13 | 152.41 | 126.78 | 78.64 | 101.14 |
| S&P 500 | 100.00 | 104.91 | 121.48 | 128.16 | 80.74 | 102.11 |

**QUESTION**: What was the value of Aimco in 2006?
**DETAIL_ANSWER**: Looking at the table, the value of Aimco in 2006 is 164.95.
**SHORT_ANSWER**: 164.95.
**(Retrieval)**

**QUESTION**: In which year did MSCI US REIT have the highest value?
**DETAIL_ANSWER**: The values for MSCI US REIT are 112.13 in 2005, 152.41 in 2006, 126.78 in 2007, 78.64 in 2008, and 101.14 in 2009. The highest value is 152.41 which is in 2006.
**SHORT_ANSWER**: 2006.
**(Data operations)**

| | 2016 | 2015 | 2014 |
|---|---|---|---|
| Fair value at grant date | $10.35 | $11.57 | $10.99 |
| Expected volatility | 21.6% | 21.0% | 25.1% |
| Dividend yield | 2.4% | 2.1% | 2.1% |
| Risk-free interest rate | 1.4% | 1.7% | 1.6% - 2.0% |
| Expected holding period, in years | 5.3 | 5.3 | 5.5 - 6.6 |

**QUESTION**: Is the fair value at grant date in 2016 higher than that in 2015?
**DETAIL_ANSWER**: The fair value at grant date in 2016 is $10.35 and in 2015 is $11.57. So, the fair value at grant date in 2016 is lower than that in 2015.
**SHORT_ANSWER**: no.
**(Selection)**

| | Year ended December 31, | | |
|---|---|---|---|
| | 2015 | 2014 | 2013 |
| | (in millions) | | |
| Total Individual Annuities(1): | | | |
| Beginning total account value | $158,664 | $154,140 | $135,342 |
| Sales | 8,780 | 10,008 | 11,513 |
| Surrenders and withdrawals | (8,415) | (8,852) | (7,727) |
| Net sales | 365 | 1,156 | 3,786 |
| Benefit payments | (1,910) | (1,799) | (1,617) |
| Net flows | (1,545) | (643) | 2,169 |
| Change in market value, interest credited and other activity | (585) | 8,666 | 19,826 |
| Policy charges | (3,589) | (3,499) | (3,197) |
| Ending total account value | $152,945 | $158,664 | $154,140 |

**QUESTION**: Please provide a detailed summary of the table.
**DETAIL_ANSWER**: In 2015, the beginning total account value was $158,664, with sales of 8,780 and surrenders and withdrawals of 8,415 resulting in net sales of 365. Benefit payments were 1,910 and net flows were -1,545. There was a change in market value, interest credited and other activity of -585 and policy charges of -3,589, leading to an ending total account value of $152,945. In 2014, similar data is presented, and in 2013 as well.
**SHORT_ANSWER**: Summary of account values, sales, etc. over 3 years.
**(Summary)**

**QUESTION**: What is the total of net sales for all three years?
**DETAIL_ANSWER**: The net sales for 2013 is 3,786, for 2014 is 1,156, and for 2015 is 365. The total is 3,786 + 1,156 + 365 = 5,307.
**SHORT_ANSWER**: 5,307.
**(Numerical calculations)**

| | Year Ended October 31, | | |
|---|---|---|---|
| | 2018 | 2017 | 2016 |
| | (in thousands, except per share price) | | |
| Shares repurchased(1) | 4,688 | 5,413 | 8,506 |
| Average purchase price per share(1) | $89.59 | $70.21 | $49.37 |
| Aggregate purchase price(1) | $420,000 | $380,000 | $420,000 |
| Reissuance of treasury stock | 3,508 | 4,404 | 4,803 |

**QUESTION**: Which year had the highest number of shares repurchased?
**DETAIL_ANSWER**: Looking at the "Shares repurchased(1)" column, we see the values 4688 for 2018, 5413 for 2017, and 8506 for 2016. Comparing these, 2016 had the highest number of shares repurchased.
**SHORT_ANSWER**: 2016.
**(Data operations)**

| Millions of Dollars | 2007 | 2006 | 2005 |
|---|---|---|---|
| Track | $1,593 | $1,487 | $1,472 |
| Capacity and commercial facilities | 534 | 510 | 509 |
| Locomotives and freight cars | 263 | 135 | 98 |
| Other | 106 | 110 | 90 |
| Total | $2,496 | $2,242 | $2,169 |

**QUESTION**: What is the total amount in 2007?
**DETAIL_ANSWER**: Looking at the table, the total amount in 2007 is $2,496.
**SHORT_ANSWER**: $2,496.
**(Retrieval)**

**QUESTION**: How much did the amount for Locomotives and freight cars increase from 2006 to 2007?
**DETAIL_ANSWER**: The amount for Locomotives and freight cars in 2006 is $135 and in 2007 is $263. The increase is $263 - $135 = $128.
**SHORT_ANSWER**: $128.
**(Numerical calculations)**

Figure 4. Visualization of SynTab sample pairs. The texts below each image represent the Q&A pairs we synthesize. Each Q&A pair includes the user-input question, detailed solution steps, and the final concise answer. The red bold text indicates the question category to which the Q&A pair belongs.

**Figure 5 (WTQ):**

Toolbar: User | SynTab-LLaVA | Qwen2-VL | GPT-4V | Table-LLaVA | TabPedia | Docowl1.5 | InternVL2

| Place | Rider | Country | Team | Points | Wins |
|---|---|---|---|---|---|
| 1 | Sylvain Geboers | Belgium | Suzuki | 3066 | 3 |
| 2 | Adolf Weil | Germany | Maico | 2331 | 2 |
| 3 | Torlief Hansen | Sweden | Husqvarna | 2052 | 0 |
| 4 | Roger De Coster | Belgium | Suzuki | 1865 | 3 |
| 5 | Joel Robert | Belgium | Suzuki | 1730 | 1 |
| 6 | Heikki Mikkola | Finland | Husqvarna | 1680 | 2 |
| 7 | Willy Bauer | Germany | Maico | 1276 | 0 |
| 8 | Gaston Rahier | Belgium | ČZ | 1112 | 0 |
| 9 | Pierre Karsmakers | Netherlands | Husqvarna | 1110 | 0 |
| 10 | Dave Bickers | United Kingdom | ČZ | 1076 | 0 |

Based on the table of '1971 Trans-AMA motocross series', answer the following question briefly. Show your answer in the JSON format {"answer": [<a list of answer strings>]}. total wins by belgian riders

**SynTab-LLaVA:** Using the table as reference, the answer is {"answer": ['7']}.

**Qwen2-VL:** {"answer": ["3", "3"]}

**GPT-4V:** I can directly count the total wins by Belgian riders based on the provided image: 1.Sylvain Geboers: 3 wins; 2.Roger De Coster: 3 wins; 3.Joel Robert: 1 win. Adding these up gives a total of 7 wins by Belgian riders. Now, I'll format this as a JSON response. {"answer": ["Total wins by Belgian riders: 7"]}

**Table-LLaVA:** According to the table, the answer is {"answer": ['10']}

**TabPedia:** 3

**Docowl1.5:** 5

**InternVL2:** {"answer": ["3"]}

**Figure 5 (FeTaQA):**

| Year | Title | Role | Channel |
|---|---|---|---|
| 2015 | Kuch Toh Hai Tere Mere Darmiyaan | Sanjana Kapoor | Star Plus |
| 2016 | Kuch Rang Pyar Ke Aise Bhi | Khushi | Sony TV |
| 2016 | Gangaa | Aashi Jhaa | &TV |
| 2017 | Iss Pyaar Ko Kya Naam Doon 3 | Meghna Narayan Vashishth | Star Plus |
| 2017–18 | Tu Aashiqui | Richa Dhanrajgir | Colors TV |
| 2019 | Laal Ishq | Pernia | &TV |
| 2019 | Vikram Betaal Ki Rahasya Gatha | Rukmani/Kashi | &TV |
| 2019 | Shaadi Ke Siyape | Dua | &TV |

Conduct table question answering task based on the given table about 'Shagun Sharma' with the table title 'Television'. What TV shows was Shagun Sharma seen in 2019?

**Ground Truth:** In 2019, Shagun Sharma played in the roles as Pernia in Laal Ishq, Vikram Betaal Ki Rahasya Gatha as Rukmani/Kashi and Shaadi Ke Siyape as Dua.

**SynTab-LLaVA:** In 2019, Shagun Sharma was seen in the TV shows Laal Ishq as Pernia, Vikram Betaal Ki Rahasya Gatha as Rukmani/Kashi and Shaadi Ke Siyaphe as Dua.

**Qwen2-VL:** Laal Ishq, Vikram Betaal Ki Rahasya Gatha, Shaadi Ke Siyape.

**GPT-4V:** In 2019, Shagun Sharma appeared in the TV shows "Laal Ishq" and "Vikram Betaal Ki Rahasya Gatha".

**Table-LLaVA:** In 2019, Shagun Sharma was seen in the role of Pernia in Laal Ishq and as Dua in Zindagi Ki Siyap.

**TabPedia:** Shagun Sharma was seen in 2 TV shows in 2019, which are 'Vikram Betaal Ki Rahasya Gatha' and 'Shaadi Ke Siyape'.

**Docowl1.5:** Colors TV &TV

**InternVL2:** Shagun Sharma was seen in the TV shows "Vikram Betaal Ki Rahasya Gatha" and "Shaadi Ke Siyape" in 2019.

Figure 5. Case study on WTQ and FeTaQA. The red rectangular boxes indicate the table regions relevant to the questions. In the WTQ, yellow highlights represent correct answers, while green highlights denote incorrect. In the FeTaQA, which requires generating a complete sentence, red text is used to indicate the keywords present in the predicted results.

**Figure 6 (TABMWP):**

Toolbar: User | SynTab-LLaVA | Qwen2-VL | GPT-4V | Table-LLaVA | TabPedia | Docowl1.5 | InternVL2

## Sit-up competition (number of sit-ups)

| Stem | Leaf |
|---|---|
| 6 | 0 1 1 3 3 5 5 6 6 9 9 |
| 7 | 2 4 5 8 |
| 8 | 0 1 6 6 6 9 |
| 9 | 0 |

Based on this table about 'Sit-up competition (number of sit-ups)', solve the following problem. Format the ultimate answer as a JSON, using the structure {"answer": "<YOUR>"}. Miss Pratt ran a sit-up competition among her P.E. students and monitored hANSWERow many sit-ups each students could do. How many people did fewer than 90 sit-ups? (Unit: people)

**SynTab-LLaVA:** Count all the leaves in the rows with stems 6, 7, and 8. You counted 21 leaves, which are blue in the stem-and-leaf plot above. 21 people did fewer than 90 sit-ups. Thus, the final answer is {"answer": "21"}.

**GPT-4V:** 1.Breakdown of the Sit-Up Data from the Table: 1) Stem 6: 11 entries (numbers 60 to 69 are represented as 0, 1, 1, 3, 3, 5, 5, 6, 6, 9, 9); 2) Stem 7: 5 entries (numbers 70 to 79 are represented as 2, 4, 5, 8); 3) Stem 8: 6 entries (numbers 80 to 89 are represented as 0, 1, 6, 6, 6, 9). 2.Total Number of People Who Did Fewer Than 90 Sit-Ups: 1) Total = Number of entries in Stem 6 + Number of entries in Stem 7 + Number of entries in Stem 8; 2)Total = 11 (Stem 6) + 5 (Stem 7) + 6 (Stem 8) = 22 people Thus, 22 people did fewer than 90 sit-ups in the competition. The final answer, formatted as requested: {"answer": "22"}

**Table-LLaVA:** Count all the leaves in the rows with stems 0, 1, 2, 3, 4, 5, 6, 7, and 8. You counted 22 leaves, which are blue in the stem-and-leaf plot above. 22 people did fewer than 90 sit-ups. Thus, the final answer is {"answer": "22"}.

**InternVL2:** {"answer": 20}

**Figure 6 (AIT-QA):**

| $ million | 2019 | 2018 | Change (%) |
|---|---|---|---|
| Order intake1 | 532.0 | 470.0 | 13.2 |
| Revenue | 503.6 | 476.9 | 5.6 |
| Gross profit | 368.6 | 344.5 | 7.0 |
| Gross margin (%) | 73.2 | 72.2 | 1.0 |
| Adjusted operating costs2 | 275.7 | 267.4 | 3.1 |
| Adjusted operating profit2 | 92.9 | 77.1 | 20.5 |
| Adjusted operating margin3 (%) | 18.4 | 16.2 | 2.2 |
| Reported operating profit | 88.6 | 57.5 | 54.1 |
| Effective tax rate4 (%) | 13.0 | 15.4 | (2.4) |
| Reported profit before tax | 89.6 | 61.2 | 46.4 |
| Adjusted basic earnings per share5 (cents) | 13.40 | 10.86 | 23.4 |
| Basic earnings per share (cents) | 12.79 | 9.14 | 39.9 |
| Free cash flow6 | 100.1 | 50.9 | 96.7 |
| Closing cash | 183.2 | 121.6 | 50.7 |
| Final dividend per share7 (cents) | 3.45 | 2.73 | 26.4 |

What was the change in closing cash? Carefully answer the question according to the table. Conclude your response with a final answer in the JSON format {"answer": [<a list of answer strings>]}.

**SynTab-LLaVA:** Based on the given table:183.2-121.6 = 61.6 Therefore, the answer is {"answer": ["61.6"]}.

**Qwen2-VL:** 50.7

**GPT-4V:** According to the table, the change in closing cash from 2018 to 2019 was 50.7%. Here is the JSON-formatted answer: {"answer": ["50.7"]}

**Table-LLaVA:** Based on the table, we need to conduct the following computations:182.2-121.6 = 60.6 Therefore, the final answer: {"answer": ["60.6"]}.

**TabPedia:** The change in closing cash was 50.7%.

**Docowl1.5:** [50.7]

**InternVL2:** The change in closing cash was 50.7%. Final answer: {"answer": ["50.7%"]}.

Figure 6. Case study on TABMWP and AIT-QA. These benchmarks involve extensive data computation tasks, requiring the outputs to present step-by-step reasoning leading to the final answer.

| 🧑 User | 🤖 SynTab-LLaVA | 🔶 Qwen2-VL | 🟣 GPT-4V | 💬 Table-LLaVA | 🤖 TabPedia | 🦉 Docowl1.5 | 🐎 InternVL2 |

| place | player | country | score | to par |
|---|---|---|---|---|
| 1 | brian watts | united states | 68 + 69 + 73 = 210 | e |
| t2 | jim furyk | united states | 70 + 70 + 72 = 212 | + 2 |
| t2 | mark o'meara | united states | 72 + 68 + 72 = 212 | + 2 |
| t2 | jesper parnevik | sweden | 68 + 72 + 72 = 212 | + 2 |
| 5 | justin rose (a) | england | 72 + 66 + 75 = 213 | + 3 |
| t6 | thomas bjärn | denmark | 68 + 71 + 76 = 215 | + 5 |
| t6 | brad faxon | united states | 67 + 74 + 74 = 215 | + 5 |
| t6 | john huston | united states | 65 + 77 + 73 = 215 | + 5 |
| t6 | tiger woods | united states | 65 + 73 + 77 = 215 | + 5 |
| t10 | david duval | united states | 70 + 71 + 75 = 216 | + 6 |
| t10 | costantino rocca | italy | 72 + 74 + 70 = 216 | + 6 |
| t10 | raymond russell | scotland | 68 + 73 + 75 = 216 | + 6 |
| t10 | katsuyoshi tomori | japan | 75 + 71 + 70 = 216 | + 6 |

Based on the table titled '1998 open championship' and the subsequent hypothesis, determine if the hypothesis is confirmed or refuted by the table. Output the final answer as JSON in the format {"answer": "<YOUR ANSWER>"}.
hypothesis: danish thomas bjørn ranked lower than tiger woods

Based on the table information, the hypothesis is not true. Thus, the answer is {"answer": "refuted"}.

{"answer": "true"}

Since both players share the same rank (T6), the hypothesis that Danish Thomas Bjørn ranked lower than Tiger Woods is refuted. Here is the answer in the requested JSON format: {"answer": "refuted"}

Based on the table information, the hypothesis is true. Thus, the answer is {"answer": "confirmed"}.

The hypothesis is refuted. Although danish thomas bjørn ranked 5th while tiger woods ranked 1st, the difference in their scores is only 4 strokes, which is smaller than the standard deviation of the scores of all the players. Therefore, it is reasonable to conclude that the difference in their ranks is due to other factors, such as the difficulty of the course or the performance of the players.

| Party | Candidate | Votes | % | ± |
|---|---|---|---|---|
| BJD | Bibhuti Bhusan Balabantaray | 55,573 | 47.85 | |
| INC | Suresh Kumar Routray | 40,648 | 35.00 | |
| Independent | Sarat Paikray | 8,756 | 7.54 | |
| BJP | Sudipta Ray | 7,812 | 6.73 | |
| SAMO | Bibhuti Bhusan Mishra | 964 | 0.83 | |
| Kalinga Sena | Pravat Pradhan | 629 | 0.54 | |
| RPD | Kirti Chandra Senapati | 610 | 0.53 | |
| BSP | Jaya Prakash Bhanja Deo | 607 | 0.52 | |
| RPI(A) | Purna Chandra Naik | 534 | 0.46 | |
| **Majority** | | 14,925 | 12.85 | |
| **Turnout** | | 1,16,142 | 60.50 | 2.81 |
| **Registered electors** | | 1,91,962 | | |
| BJD hold | | | | |

Examine the following Wikipedia table, which is titled 'Jatani (Odisha Vidhan Sabha constituency)' and its section is titled '2009 Election Results'.
Some important table cells are highlighted with a different background color. Please create a single sentence that summarizes the highlighted table cells.

In 2009 election, Biju Janata Dal candidate Bibhuti Bhusan Balabantaray defeated Indian National Congress candidate Suresh Kumar Routray by a margin of 14,925 votes.
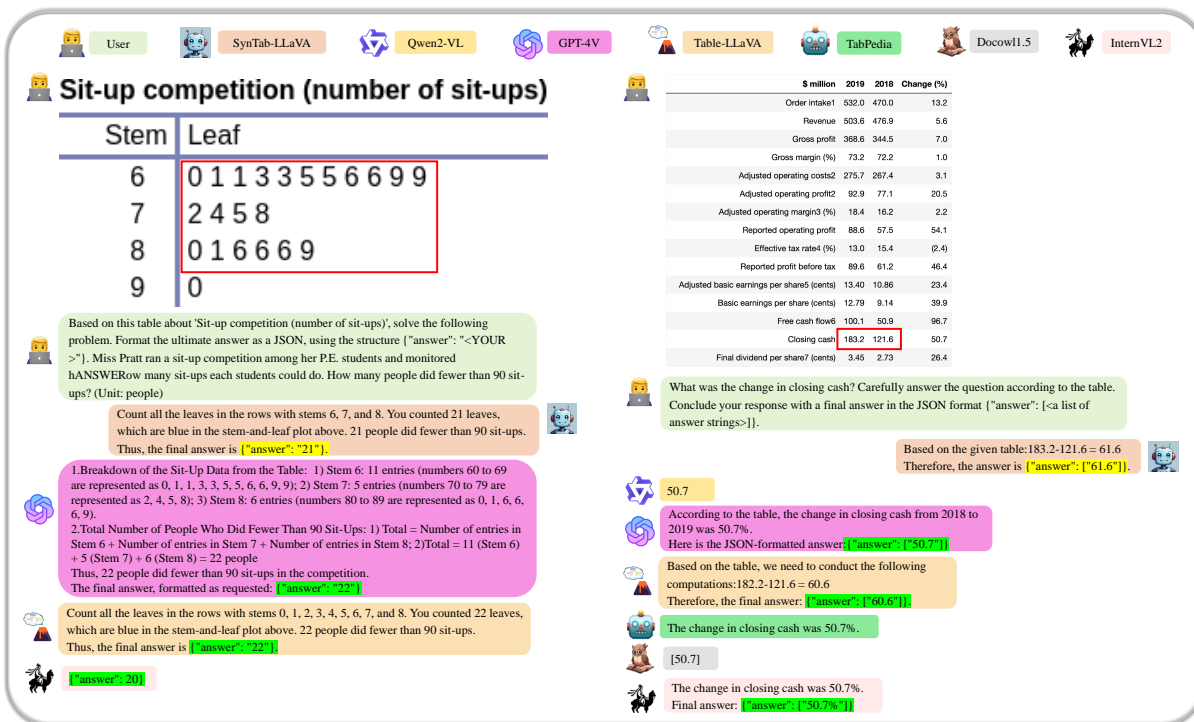
Biju Janata Dal candidate Bibhuti Bhusan Balabantaray won the Jatani Vidhan Sabha constituency with 55,573 votes.

In the 2009 election results for the Jatani (Odisha Vidhan Sabha constituency), Bibhuti Bhusan Balabantaray of the BJD party received 55,573 votes, securing 47.85% of the votes, and the majority needed to win was 14,925 votes.

In the 2009 election, Biju Janata Dal candidate Bibhuti Bhusan Balabantaryy defeated Indian National Congress candidate Suresh Kumar Routray by a margin of 14,925 votes.

The highlighted table cells show the results of the 2009 Jatani (Odisha Vidhan Sabha constituency) election, with the BJD candidate Bibhuti Bhusan Balabantaray winning the majority of votes and the overall turnout being 60.50%.
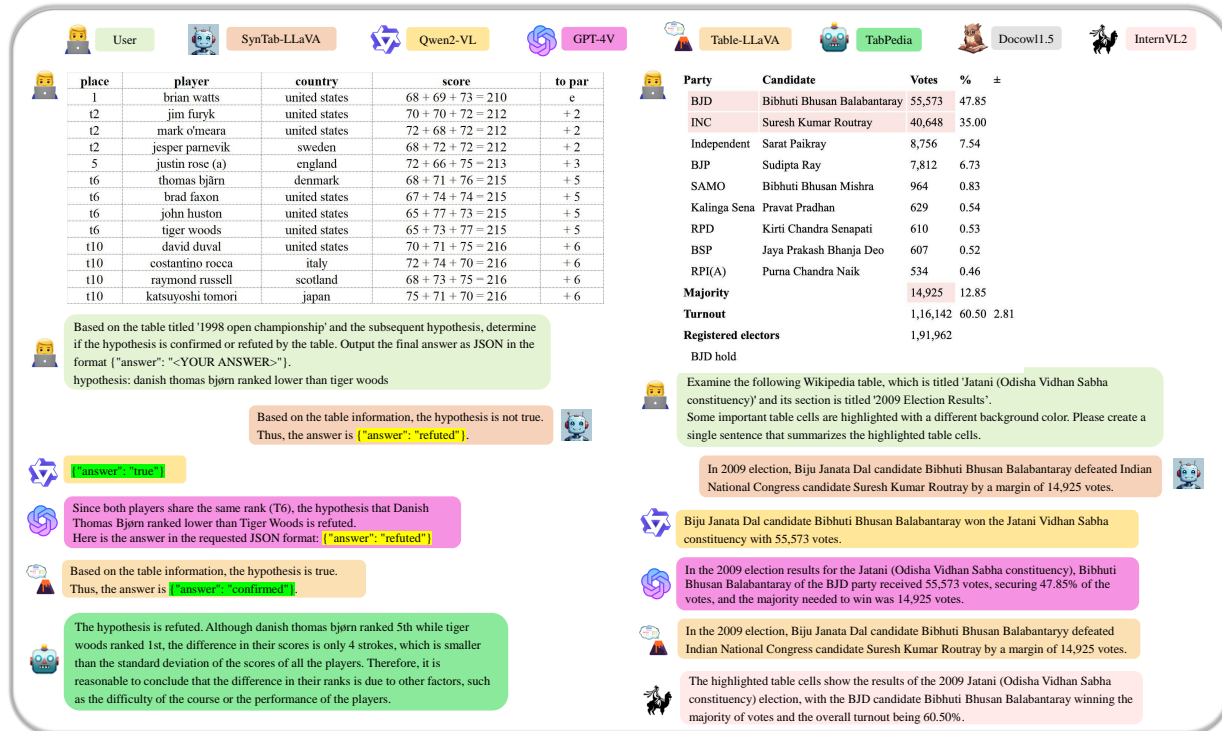
Figure 7. Case study on TabFact and ToTTo. The ToTTo benchmark not only requires the model to perceive color information in images but also demands strong summarization capabilities, posing a challenge to the captioning abilities of existing MLLMs.



| 🧑 User | 🤖 SynTab-LLaVA | 🔶 Qwen2-VL | 🟣 GPT-4V | 💬 Table-LLaVA | 🤖 TabPedia | 🦉 Docowl1.5 | 🐎 InternVL2 |

| Year | Rider | Victories | Bike | Manufacturer's Championship |
|---|---|---|---|---|
| 1994 | Carl Fogarty | 11 | Ducati 916 | Ducati |
| 1995 | Carl Fogarty | 13 | Ducati 916 | Ducati |
| 1996 | Troy Corser | 7 | Ducati 916 | Ducati |
| 1998 | Carl Fogarty | 3 | Ducati 916 | Ducati |
| 1999 | Carl Fogarty | 11 | Ducati 996 | Ducati |
| 2000 | *(Colin Edwards)* | *(7)* | *(Honda RC51)* | Ducati |
| 2001 | Troy Bayliss | 6 | Ducati 996 | Ducati |

Please determine the total count of rows and columns in the provided table, respectively. Provide the final answer in the JSON structure, using the format {"row_number": "m", "column_number": "n"}.

This table has 8 rows and 5 columns. Thus, the final answer is {"row_number": "8", "column_number": "5"}.

{"row_number": "7", "column_number": "4"}

The provided table contains a total of 7 rows and 4 columns. Here is the final answer in the requested JSON structure: {"row_number": "7", "column_number": "4"}

This table has 8 rows and 3 columns. Thus, the final answer is {"row_number": "8", "column_number": "3"}.

The table has 7 rows and 5 columns.

[{"row_number": 11, "column_number": "victories"}, {"row_number": 13, "column_number": "victories"}, {"row_number": 7, "column_number": "victories"}, {"row_number": 3, "column_number": "victories"}]

To determine the total count of rows and columns in the provided table, we need to count each row and each column separately. Counting Rows: 1. 1994 2. 1995 3. 1996 4. 1998 5. 1999 6. 2000 7. 2001. There are 7 rows in total. Counting Columns: 1. Year 2. Rider 3. Victories 4. Bike 5. Manufacturer's Championship. There are 5 columns in total. The final answer in JSON format is: {"row_number": "7", "column_number": "5"}

| | linked cctip tumour cohort | tumours with one or more surgical interventions | | |
|---|---|---|---|---|
| | number | number | percentage | |
| canada | 44,970 | 37,315 | 83.0 | |
| newfoundland and labrador | 1,495 | 1,295 | 86.5 | |
| prince edward island | 300 | 250 | 82.5 | |
| nova scotia | 2,330 | 1,930 | 82.9 | |
| new brunswick | 1,545 | 1,275 | 82.5 | |
| ontario | 21,510 | 17,440 | 81.1 | |
| manitoba | 2,385 | 1,945 | 81.7 | |
| saskatchewan | 2,070 | 1,745 | 84.2 | |
| alberta | 5,145 | 4,320 | 84.0 | |
| british columbia | 8,050 | 6,995 | 86.9 | |
| yukon | 45 | 40 | 87.0 | |
| northwest territories | 70 | 60 | 84.5 | |
| nunavut | 25 | 20 | 88.5 | |

Here's what I need you to do:
In the image of this table, identify all the merged cells. Provide the coordinates of these merged regions in the JSON {'top-left':(x, y), 'bottom-right':(z, w)}, with '(x, y)' and '(z, w)' indicating the (row_id, column_id) of the top-left and bottom-right cells, respectively. Remember, row and column numbering starts at 1.

The merged cell regions are as follows: {'top-left':(1, 1), 'bottom-right':(2, 1)}; {'top-left':(1, 3), 'bottom-right':(1, 4)}.

{"top-left":{"row_id":1,"column_id":1},"bottom-right":{"row_id":1,"column_id":2}}

…This indicates that the cell containing "Canada" spans across all four header columns in the first row. Here is the JSON output with the merged region coordinates: { "top-left": [1, 1], "bottom-right": [1, 4] }

This table does not contain any merged cells.

Here is a list of all the merged cells in the table: [0.000,0.402,0.243,0.422]

{'top-left':(1, 1), 'bottom-right':(17, 17)}

To identify all the merged cells in the table, we need to look for cells that are merged horizontally or vertically. Merged cells in a table are those that are combined into a single cell, which can be identified by the presence of a merged cell icon in the table editor... let's convert these merged cells into the required JSON format:json{ "merged_cells":
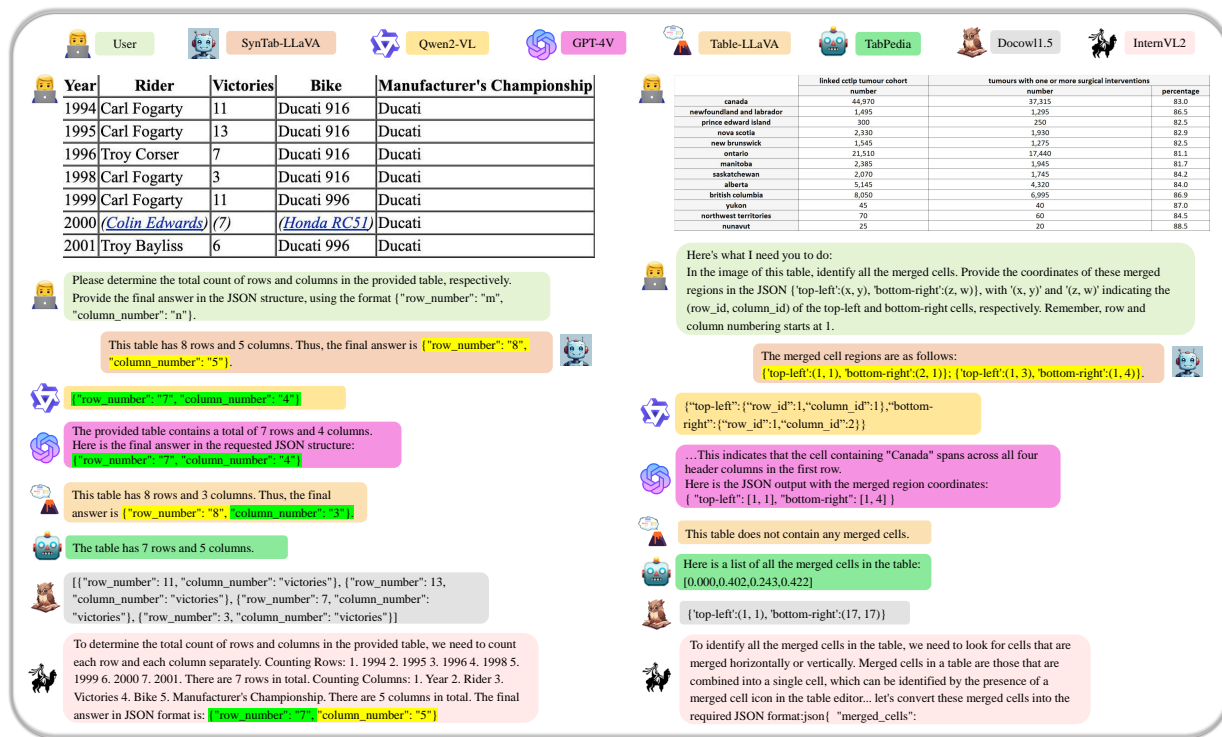
Figure 8. Case study on TSD and MCD. These two benchmarks primarily evaluate a model's ability to perceive the structure of table images. The results indicate that existing MLLMs perform poorly on such fundamental perception tasks.
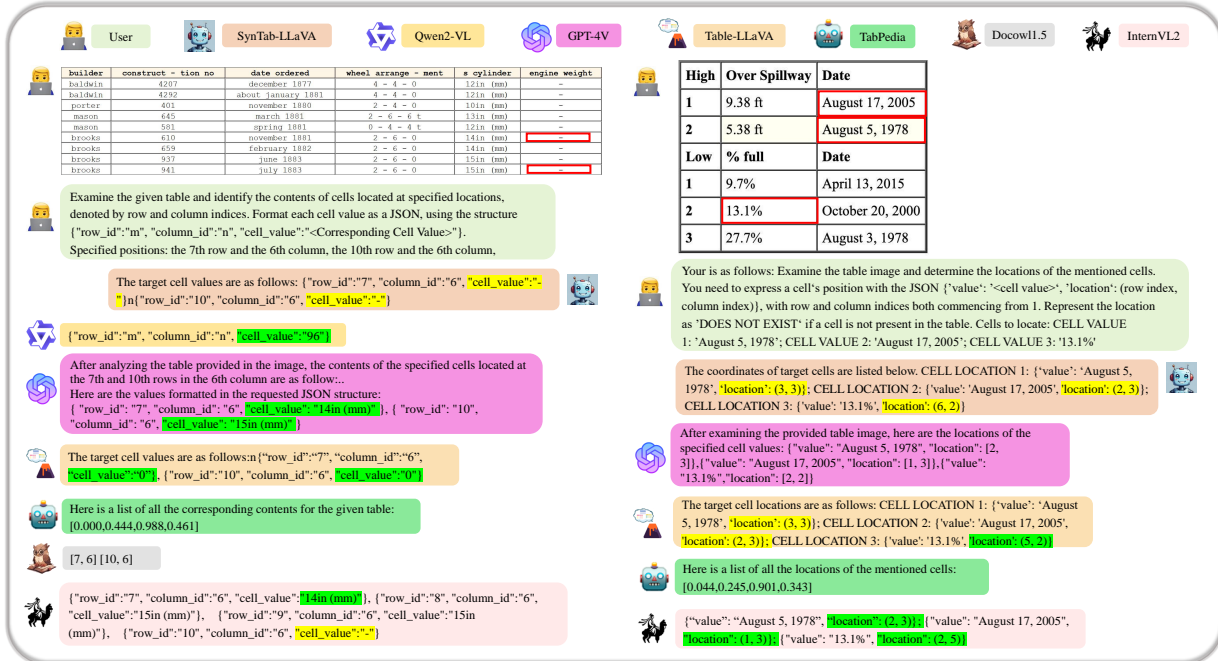
Figure 9. Case study on TCE and TCL. The visualization results above demonstrate that existing MLLMs exhibit significant gaps compared to our method in locating and recognizing table cell content.
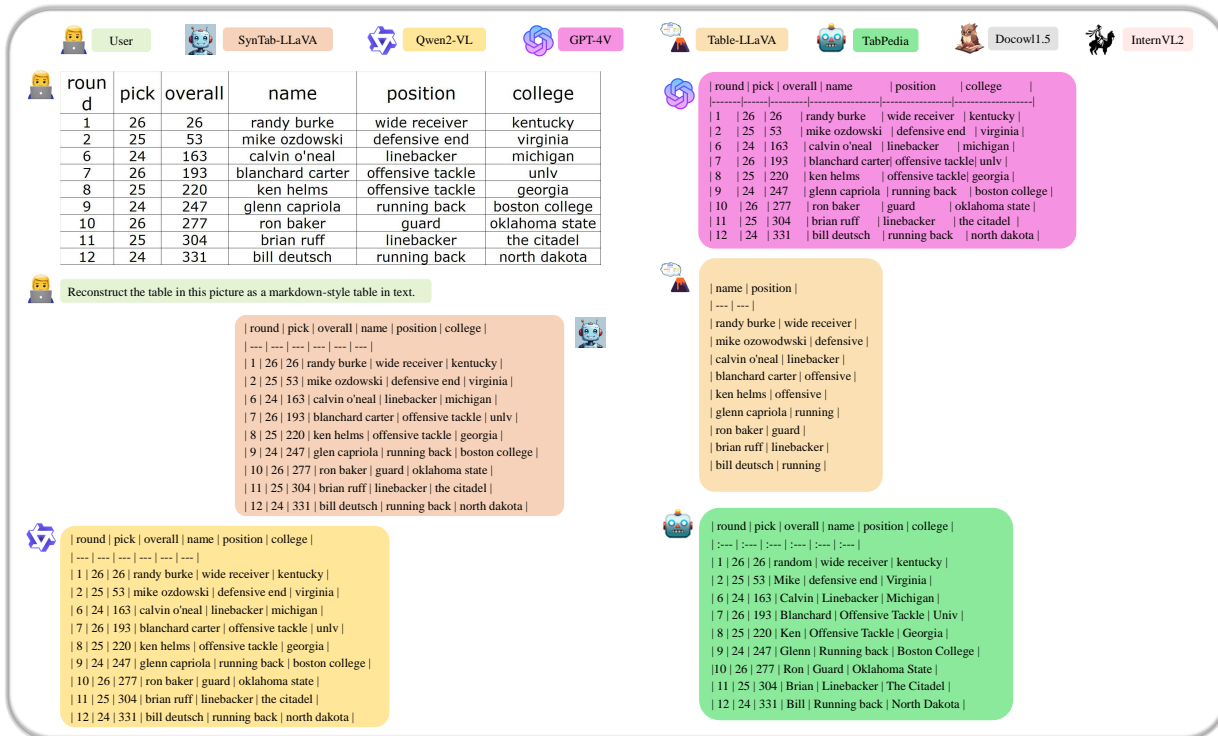
Figure 10. Case study on TR. The general MLLMs demonstrate strong performance in converting images to Markdown, primarily because they are pre-trained on large-scale Markdown datasets. Compared to the tabular MLLMs Table-LLaVA and TabPedia, our approach achieves a significant performance improvement.

numerical values; multiplication: calculating the product of numerical values; division: determining the quotient of numerical value; simple combinations: performing basic operations involving combinations of addition, subtraction, multiplication, and division.

**Complex calculations** involve operations that combine multiple tasks, including: data retrieval, sorting, filtering, mathematical operations, and so on. Complex calculations integrate these tasks to perform more intricate data analysis.

### B.4. Free Answering

**Free table question answering.** This subcategory requires retrieval and reasoning over multiple sources of information in table, and the answer should integrate both facts and inferences into a coherent sentence that answers the question.

### B.5. Selection

**Multiple choice** in the context of MTU involves selecting the correct answer from a list of given options based on the data presented in the table.

**Table fact verification** involves checking the accuracy and validity of information presented in a table against known facts or expected results. The main aspects of table fact verification include: cross-checking, consistency check, error detection, logical verification, *etc*.

### B.6. Summary

**Table Summary** involves providing a concise and coherent description of the key information presented in a table.

## C. Visualization

### C.1. Rendered Table Images

We visualize the synthetic table images generated in this paper. As shown in Fig. 3, our rendered images exhibit diverse characteristics, including different font styles, varying highlighted regions, and a variety of table visual layout, such as only row separators, only column separators, both row and column separators, and row segmentation using background colors, *etc*. Compared to existing MTU images, our synthesized images are even more diverse and better match common table images found in real-world scenarios. Additionally, our images also include complex hierarchical tables, as seen in the last row of the figure. These hierarchical table images facilitate the model's learning of complex table structures, thereby enhancing its capability to understand tables in complex scenarios. From the visualizations above, we conclude that, through extensive data augmentation, the synthesized images closely resemble real-world table images, thus validating the effectiveness and authenticity of the synthetic data proposed in this paper.

### C.2. Synthesized SynTab Samples

As shown in Fig. 4, we visualize some sample pairs from SynTab. It can be observed that the questions generated for each image are not fixed but are randomly selected based on probability to determine the type of question to generate. This design prevents the dataset from being dominated by specific question types, enriches the diversity of samples, and enhances the generalization and robustness of subsequent MLLMs in table understanding tasks. Additionally, our generated answers include a detailed answer that provides step-by-step solutions to the questions. Such Q&A pairs enable the model to learn how to approach table understanding tasks incrementally rather than producing a direct answer, thereby avoiding potential performance degradation caused by the lack of intermediate reasoning steps.

### C.3. Comparison with Other MLLMs

To compare our method with existing MLLMs in table understanding tasks, we conduct visualized experiments across multiple benchmark datasets, with the results shown in Fig. 5-10. The experimental results demonstrate that our model SynTab-LLaVA achieves significant performance improvements across various benchmarks when compared to Qwen2-VL, GPT-4V, Docowl1.5, InternVL2, as well as tabular MLLMs such as Table-LLaVA and TabPedia. For example, in Fig. 5, the input in the left sample requires the model to retrieve all the "wins" cells corresponding to riders from Belgium in the table, then sum these "wins" to generate the answer. It can be observed that, apart from GPT-4V and our SynTab-LLaVA, the other models make incorrect predictions. A closer analysis of GPT-4V's answer shows that while it predicts the correct result, it misses a rider, Gaston Rahier, whose "wins" are 0, making its answer unsatisfactory. As for the FeTaQA, We use red text to highlight the keywords in the predicted answer sentence that match the ground truth. As shown in Fig. 5, our method successfully matches the most keywords, indicating that our model not only generates concise answers but also produces coherent sentences capable of handling a variety of MTU tasks. These qualitative analysis validate the effectiveness of the proposed method and highlights its potential contribution to the multimodal community by providing a new MTU dataset to enhance model performance in multimodal table understanding.