

## A. Fine-tuning Datasets

During our fine-tuning stage, VLMs are optimized using a variety of datasets with different tasks. These datasets, introduced in Table 2, include InfographicVQA [25], Kleister Charity [31], WikiTableQuestions [28], VizWiz-VQA [10], and ST-VQA [4], and are briefly described as follows:

**InfographicVQA (InfoVQA) [25]:** This dataset is a collection of over five thousand infographic images, along with a large number of question-answer pairs. These infographics are sourced from various web domains and feature diverse layouts and designs. The InfographicVQA challenges vision language models to interpret and reason over complex visual documents, often necessitating understanding of graphical elements, data visualization, reasoning, and arithmetic skills.

**Kleister Charity (KLC) [31]:** This dataset consists of annual financial reports from UK charity organizations. The task involves key information extraction (KIE) such as charity names, addresses, charity numbers, and reporting dates. Primarily comprising scanned documents, this dataset poses challenges due to its length, diverse layout, and the necessity to interpret both text and structural features.

**WikiTableQuestions (WTQ) [28]:** This dataset includes question and answer pairs collected from thousands of HTML tables extracted from Wikipedia. The questions are designed to be complex, requiring multi-step reasoning and various data operations such as comparison, aggregation, and arithmetic computation.

**VizWiz-VQA (VizWiz) [10]:** Comprising a large number of question-answer pairs, this dataset features images captured by blind individuals using mobile phones and spoken questions about those images. Unique in terms of its image quality, which is often blurred, and the nature of its questions, a significant portion of the images are unanswerable due to poor image quality.

**ST-VQA [4]:** Designed specifically for understanding textual information within natural images, this dataset requires models to read and interpret scene text to accurately answer questions. It includes a large collection of images sourced from various public datasets such as COCO-text, Visual Genome, and ImageNet, challenging models to comprehend images across a wide range of scenarios and textual appearances within images.

## B. More Qualitative Results

Figure 5 presents additional quantitative results derived from various evaluation datasets comparing QA-ViT method [9] and our QID method, implemented with Qwen-VL-Chat [3]. These results highlight the effectiveness of our method in enhancing the vision model’s ability to identify relevant visual cues and improve comprehension in both text-rich and natural scene environments.

Furthermore, we outline the limitations of our approach in Figure 6. Although our method aids in enhancing understanding of text-rich images, it does not significantly improve the model’s reasoning and arithmetic capabilities. Consequently, our future research will focus on refining the model’s ability to perform complex reasoning tasks more effectively in dense-text settings.

## C. Broader Impact

The enhanced capabilities of vision-language models (VLMs) offer substantial promise for improving document comprehension in environments with extensive textual content. However, the interaction between question embeddings and vision representations remains relatively unexplored. Our approach encourages this interaction with limited fine-tuning samples while preserving the structural integrity of pre-trained VLMs. It also minimizes the necessity for extensive retraining, thereby reducing the computational resources required for deploying sophisticated AI solutions. Additionally, our method’s efficiency with limited data can decrease both the time and costs involved in annotating large datasets, enhancing the accessibility and affordability of advanced document understanding technologies. We encourage the research community to further explore and adopt our QID for text-intensive tasks, anticipating significant benefits in various applications.

[illegible]

Record	Finish	Manager	Playoffs
68-66	6th	Torrey Poffo / Bill Wolfe	none
77-76	5th	Jim Norren	none
80-79	4th	Bill Wolfe	none
90-98	15th	<b>Bob Lemon</b>	lost in 1st
75-72	6th (H)	George Case	lost in 1st
118-84	12th	George Case	lost in 1st
60-81	12th (H)	Wayne Stanger	lost in 1st
74-69	1st	Ed Baller	lost in 1st
78-87	3rd	Chuck Tanner	lost in 1st
98-48	1st	Chuck Tanner	lost in League
73-79	11th	Bill Wolfe	lost in 1st
74-74	5th	Ricky Bickel	lost in 1st
72-74	5th	Ricky Bickel / Warren Hacker / Ray Harshbarger	lost in 1st
87-77	4th	Ray Harshbarger	lost in 1st
89-88	4th	Ray Harshbarger	lost in 1st
86-84	4th	Ray Harshbarger	lost in 1st
74-67	2nd	Dick Phillips	lost in League
56-82	8th	Doug Decker	lost in League
74-75	8th	Dick Phillips	lost in League
74-65	9th	Doug Decker	lost in League
72-65	10th (H)	Doug Decker	lost in 1st
74-71	10th	Doug Decker	lost in 1st
72-71	5th	Tom Trebbelheim	lost in 1st
73-73	1st	Tom Trebbelheim	lost in 1st
84-59	1st	Torrey Poffo	lost in 1st
65-79	7th	Torrey Poffo	lost in 1st

Week	Date	Opponent	Result	Attendance
2	September 13, 2007	San Diego Chargers	W 20-13	56,438
3	September 20, 2007	at Seattle Seahawks	L 41-14	63,667
4	September 27, 2007	Minnesota Vikings	cancelled	
5	October 4, 2007	at Los Angeles Raiders	L 20-17	55,868
6	October 11, 2007	at Miami Dolphins	L 42-6	10,768
7	October 18, 2007	Dallas Cowboys	L 26-13	20,496
8	October 25, 2007	at San Diego Chargers	L 42-20	47,972
9	November 1, 2007	at Chicago Bears	L 31-28	63,498
10	November 8, 2007	at Pittsburgh Steelers	L 27-24	68,411
11	November 15, 2007	New York Jets	L 29-13	44,611
12	November 22, 2007	Green Bay Packers	L 23-3	38,418
13	November 29, 2007	at San Diego	L 20-17	55,868
14	December 6, 2007	at Cincinnati Bengals	L 27-13	45,403
15	December 13, 2007	at New York Jets	W 16-10	62,834
16	December 19, 2007	at Denver Broncos	L 20-17	76,620
17	December 27, 2007	Seattle Seahawks	W 40-20	20,370

[illegible][illegible]

Rank	Cyclist	Team	Time	100 Points
1	Alphonse Valenciennes (BEL)	Concorde d'Esperance	2h 25' 20"	47
2	Alphonse Bollenbacher (BEL)	Team CSC-Saxo Bank	31	29
3	Emilio Bianchi (ITA)	Comandante	31	29
4	Francis Bitterli (ITA)	Quick Step	31	29
5	Francis Bollenbacher (ITA)	Esperance	31	29
6	Francis Bollenbacher (ITA)	Esperance	31	29
7	Samuel Sanchez (ITA)	Esperance	31	29
8	William Rancoule (ITA)	Esperance	31	29
9	Francis Bollenbacher (ITA)	Esperance	31	29
10	Francis Bollenbacher (ITA)	Esperance	31	29

**Question:** What was the total number of points won by Franco Pellizzotti?

**Owen-VL-Chart (QA-VIT):** 20

**Owen-VL-Chart (QID):** 15

Year	Date	Opponent	Place (score)	Time (score)	Notes	Attendance
1971	12/12	Winnipeg	W 1-0	1:00	First game	2,000
1972	1/13	Winnipeg	W 2-0	1:00		2,000
1973	1/13	Winnipeg	W 2-0	1:00		2,000
1974	1/13	Winnipeg	W 2-0	1:00		2,000
1975	1/13	Winnipeg	W 2-0	1:00		2,000
1976	1/13	Winnipeg	W 2-0	1:00		2,000
1977	1/13	Winnipeg	W 2-0	1:00		2,000
1978	1/13	Winnipeg	W 2-0	1:00		2,000
1979	1/13	Winnipeg	W 2-0	1:00		2,000
1980	1/13	Winnipeg	W 2-0	1:00		2,000
1981	1/13	Winnipeg	W 2-0	1:00		2,000
1982	1/13	Winnipeg	W 2-0	1:00		2,000
1983	1/13	Winnipeg	W 2-0	1:00		2,000
1984	1/13	Winnipeg	W 2-0	1:00		2,000
1985	1/13	Winnipeg	W 2-0	1:00		2,000
1986	1/13	Winnipeg	W 2-0	1:00		2,000
1987	1/13	Winnipeg	W 2-0	1:00		2,000
1988	1/13	Winnipeg	W 2-0	1:00		2,000
1989	1/13	Winnipeg	W 2-0	1:00		2,000
1990	1/13	Winnipeg	W 2-0	1:00		2,000
1991	1/13	Winnipeg	W 2-0	1:00		2,000
1992	1/13	Winnipeg	W 2-0	1:00		2,000
1993	1/13	Winnipeg	W 2-0	1:00		2,000
1994	1/13	Winnipeg	W 2-0	1:00		2,000
1995	1/13	Winnipeg	W 2-0	1:00		2,000
1996	1/13	Winnipeg	W 2-0	1:00		2,000
1997	1/13	Winnipeg	W 2-0	1:00		2,000
1998	1/13	Winnipeg	W 2-0	1:00		2,000
1999	1/13	Winnipeg	W 2-0	1:00		2,000
2000	1/13	Winnipeg	W 2-0	1:00		2,000
2001	1/13	Winnipeg	W 2-0	1:00		2,000
2002	1/13	Winnipeg	W 2-0	1:00		2,000
2003	1/13	Winnipeg	W 2-0	1:00		2,000
2004	1/13	Winnipeg	W 2-0	1:00		2,000
2005	1/13	Winnipeg	W 2-0	1:00		2,000
2006	1/13	Winnipeg	W 2-0	1:00		2,000
2007	1/13	Winnipeg	W 2-0	1:00		2,000
2008	1/13	Winnipeg	W 2-0	1:00		2,000
2009	1/13	Winnipeg	W 2-0	1:00		2,000
2010	1/13	Winnipeg	W 2-0	1:00		2,000
2011	1/13	Winnipeg	W 2-0	1:00		2,000
2012	1/13	Winnipeg	W 2-0	1:00		2,000
2013	1/13	Winnipeg	W 2-0	1:00		2,000
2014	1/13	Winnipeg	W 2-0	1:00		2,000
2015	1/13	Winnipeg	W 2-0	1:00		2,000
2016	1/13	Winnipeg	W 2-0	1:00		2,000
2017	1/13	Winnipeg	W 2-0	1:00		2,000
2018	1/13	Winnipeg	W 2-0	1:00		2,000
2019	1/13	Winnipeg	W 2-0	1:00		2,000
2020	1/13	Winnipeg	W 2-0	1:00		2,000
2021	1/13	Winnipeg	W 2-0	1:00		2,000
2022	1/13	Winnipeg	W 2-0	1:00		2,000

**Question:** What date was the only game played on soldier field?

**Qwen-VL-Chat (QA-VT):** 10 October 27

**Qwen-VL-Chat (QID):** October 23

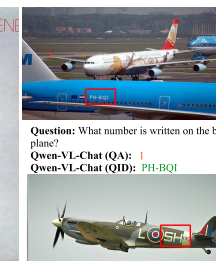
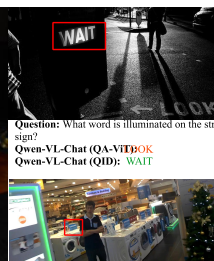
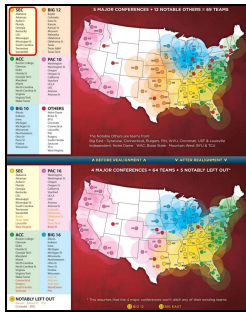
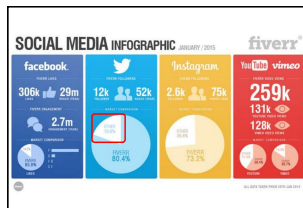


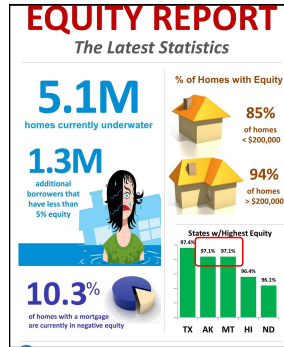
Figure 5. More qualitative results between between QA-ViT and our QID with Qwen-VL-Chat model. Image regions with answers are highlighted.



**Question:** How many teams did SEC conference have before re-alignment?  
**Qwen-VL-Chat (QID):** 14  
**Ground Truth:** 12



**Question:** What is the percentage market share others have in comparison to Fiverr in Twitter?  
**Qwen-VL-Chat (QID):** 80.4%  
**Ground Truth:** 19.6%



**Question:** How many states have 97.1% equity?  
**Qwen-VL-Chat (QID):** 3  
**Ground Truth:** 2

Iteration	Year	Date	Location	Theme
1st	1972 <sup>[20]</sup>	6 May-20 May	Suva, Fiji	"Preserving culture"
2nd	1976 <sup>[13]</sup>	6 March-13 March	Rotorua, New Zealand	"Sharing culture"
2nd	1980 <sup>[21]</sup>	30 June-12 July	Port Moresby, Papua New Guinea	"Pacific awareness"
4th	1985 <sup>[32]</sup>	29 June-13 July	Tahiti, French Polynesia	"My Pacific"
5th	1988 <sup>[34]</sup>	14 August-24 August	Townsville, Australia	"Cultural interchange"
6th	1992 <sup>[35]</sup>	16 October-27 October	Rarotonga, Cook Islands	"Seafaring heritage" <sup>[36]</sup>
7th	1996 <sup>[37]</sup>	6 September-23 September	Aloa, Samoa	"Unveiling treasures"
8th	2000 <sup>[38]</sup>	23 October-3 November	Nouméa, New Caledonia	"Words of past, present, future" <sup>[39]</sup>
9th	2004 <sup>[39]</sup>	22 July-31 July	Koror, Palau	"Warfare, Regeneration, Celebration" <sup>[23]</sup>
10th	2006 <sup>[22]</sup>	26 July-2 August	Papa Paga, American Samoa	"Threaded the Gossamer Web"
11th	2012	1-14 July	Honiara, Solomon Islands	"Culture in Harmony with Nature" <sup>[23]</sup>
12th	2016	TBA	Tamari, Guam	"TBA"
13th	2020 <sup>[24]</sup>	TBA	TBA, Hawaii	"TBA"

**Question:** What is the average length of the festival as of 2012?  
**Qwen-VL-Chat (QID):** 10 years  
**Ground Truth:** 14 days

Figure 6. Failure cases of QID on documents and questions require arithmetic and reasoning skills. Image regions with answers are highlighted.