
Heterogeneous Graph Network for Multi-page Document VQA

Team: Ting-Chih and Hanwen

Agenda

- Motivation
- MP-DocVQA Dataset
- Competition task baseline
- Methods
- Results

Single-page DocVQA

Inputs

Question

What is the idea behind the consumer relations efficiency team?

July 31, 2003

Consumer Relations Efficiency Team

Team:

Nancy Bonland	Maranda McManis	Vivian Kohnen
Cary Hilda	Nancy McManis	Yvonne Wilford
Todd Hultmark	Brian O'Brien	
Chris Hansen	Doreen Wallick	

Objective:

Balance cost efficiency with quality customer service

Our Process:

• Essential Customer Relations core services:

- | | |
|----------------------|---------------------|
| - Warm Interactive | - AT&T |
| - Calling - Live Rep | - Brand Chargebacks |
| - YA - Live Rep | - Verio |
| - Customers | |

Response for Customer Contact:	%
Request Catalog/Order Form	100%
Order Status	100%
Making List Request (add, change, update add, cancel change)	100%
Transfer of Questions (bill, add'l service questions, forms, Web issues, etc.)	25%
Product Issues	40%
Reconnect DMS, connections, full circuit	10%
CLP Issues - card request, stream in issues	15%
Other	10%

Document Question Answering Model

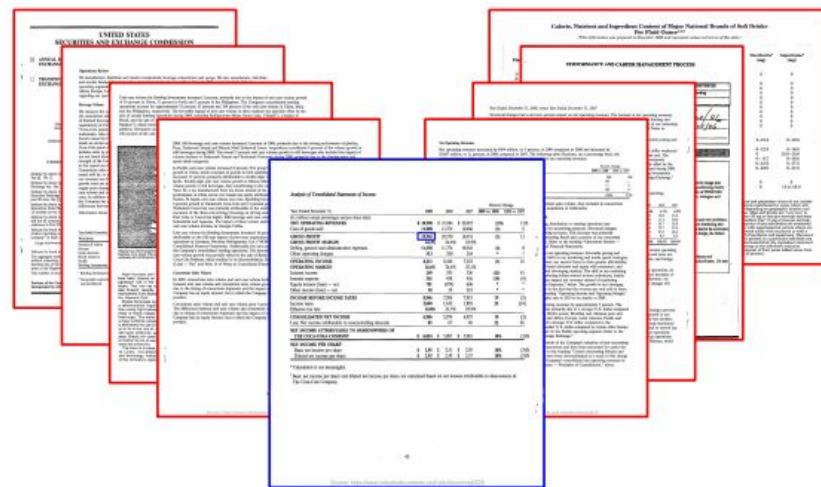
Output

Answer

Balance cost efficiency with quality customer service

Motivation

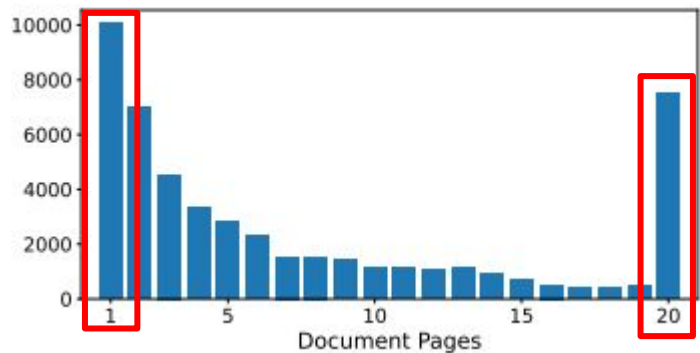
- Robust Reading Competition
- Existing work on DocVQA only considers single-page document
- In real scenes, DocVQA applies to banking, insurance and public administration
- To address the multi-page DocVQA, we should learn the reasoning from handwritten text, layout and graphical elements



Q: What was the gross profit in the year 2009?
A: \$19,902

MP-DocVQA Dataset

- MP-DocVQA comprises 46K questions posed over 48K images of scanned pages that belong to 6K industry documents
- Dataset(75/12.5/12.5):
 - Training: 36,230
 - Validation: 5,187
 - Testing: 5,019



Dataset	Questions	Documents	Pages (Images)	Avg. pages per question	Question Avg. length	Answer Avg. length	Document Avg. OCR Tokens
SingleDocVQA [16]	50K	6K	12K	1.00	9.49	2.43	151.46
VisualMRC [22]	30K	10K	10K	1.00	10.55	9.55	182.75
InfographicsVQA [15]	30K	5.4K	5.4K	1.00	11.54	1.60	217.89
DuReaderVis [19]	15K	158K	158K	1.3K	9.87	180.54	1968.21
DocCVQA [23]	20	14K	14K	14K	14.00	12.75	509.06
TAT-DQA [31]	16K	2.7K	3K	1.07	12.54	3.44	550.27
MP-DocVQA (ours)	46K	6K	48K	8.27	9.90	2.20	2026.59

MP-DocVQA Dataset example

11:14 to
11:39 a.m. Coffee Break
Coffee will be served for men and
women in the lobby adjacent to
exhibit area. Please move into
exhibit area. **(Exhibits Open)**

11:39 a.m. TRRF GENERAL SESSION (PART I)
Presiding: Lee A. Waller
TRRF Vice President

11:39 to
11:44 a.m. "Introductory Remarks"
Lee A. Waller, TRRF Vice Presi-
dent

11:44 a.m. Individual Interviews with TRRF
to Public Board Members and Sci-
12:25 p.m. entific Advisory Council Mem-
bers

Conducted by TRRF Treasurer
Philip G. Kuehn to get answers
which the public refrigerated
warehousing industry is looking
for. Plus questions from the floor.
Dr. Emil M. Mrak, University of Cal-
ifornia, Chairman, TRRF Board;
Sam R. Cecil, University of Georgia
College of Agriculture; Dr. Stanley
Charm, Tufts University School of
Medicine; Dr. Robert H. Cotton, ITT
Continental Baking Company; Dr.
Owen Fennema, University of Wis-
consin; Dr. Robert E. Hardenburg,
USDA.

12:25 to
12:58 p.m. Questions and Answers

12:58 to
4:00 p.m. Exhibits Open
Capt. Jack Stoney Room

2:00 to
5:00 p.m. TRRF Scientific Advisory
Council Meeting
Ballroom Foyer

Source: <https://www.industrydocuments.ucsf.edu/docs/hxnp0227>

MONDAY, MAY 15

8:15 to
8:56 a.m. Exhibits Open
Capt. Jack Stoney Room

8:58 a.m. OPENING GENERAL SESSION
Learnington Hall
(Ladies are invited to hear Dr.
Klaus and Dr. Feinberg)
Presiding: Charles D. Nesbit,
IARW Chairman

8:58 to
9:03 a.m. "Opening Remarks"
Charles D. Nesbit, IARW Chair-
man

9:03 to
9:07 a.m. Report of IARW Nominating Com-
mittee
James G. Talbot, Chairman

Report of TRRF Nominating Com-
mittee
Willis S. McLeese, Chairman

9:08 to
9:53 a.m. "Be Tomorrow's Person Today"
Dr. Gunther Klaus, Managing
Director, Institute for Advanced
Planning, Beverly Hills, Califor-
nia

9:53 to
10:08 a.m. Questions and Answers

10:09 to
10:59 a.m. "People Are Your Future. For
Good or Ill, You and Your Com-
pany Depend on Their Wisdom,
Their Motivation and Their
Energy"
Dr. Mortimer R. Feinberg, Chair-
man of the Board, BFS Psycho-
logical Associates, Inc., New
York City

10:59 to
11:14 a.m. Questions and Answers

Source: <https://www.industrydocuments.ucsf.edu/docs/hxnp0227>

PROGRAM

SUNDAY, MAY 14

8:30 a.m. IARW Board of Directors Meeting
Coggins Point

10:00 a.m. Registration Desk Opens
Lower Lobby

Noon to
4:00 p.m. Exhibits Open
Capt. Jack Stoney Room

Noon IARW Board of Directors
TRRF Board of Governors
TRRF Scientific Advisory Council

Noon—Reception
12:30—Luncheon
Learnington Hall North

1:00 to
4:00 p.m. Ladies' Hospitality Suite
Possum Point

1:30 p.m. TRRF Board of Governors
TRRF Scientific Advisory
Council
Joint Meeting
Learnington Hall South

6:45 to
8:00 p.m. Welcoming Reception
Pool Deck

Co-sponsored by IARW and
the following Associate Members:
American Isowall Corporation
American Refrigeration
Contractors
Clark Door Company
C. T. Hogan & Company
Kramer Trenton Company
Pittsburgh Corning Corporation
Stuart V. Smith Company
Superior Industries, Inc.

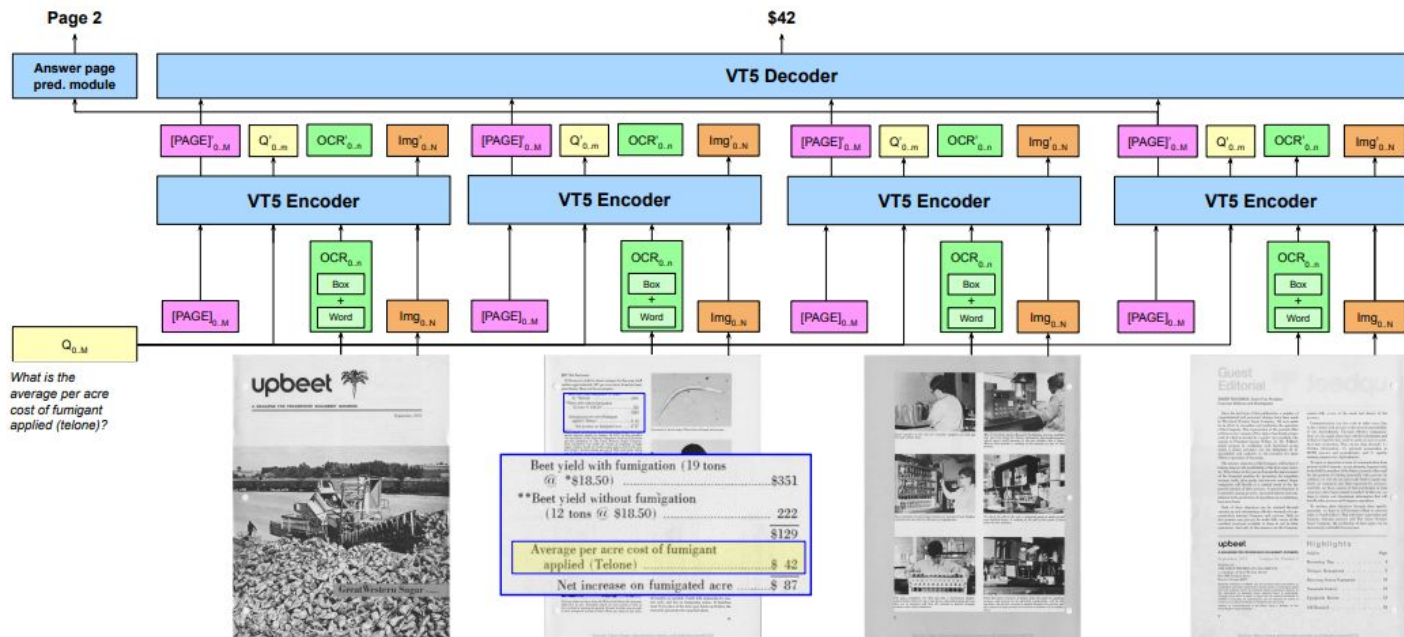
Source: <https://www.industrydocuments.ucsf.edu/docs/hxnp0227>

Q: What time is the coffee break?

A: 11.14 to 11.39 a.m.

Competition task baseline: Hi-VT5

- Hierarchical Visual T5



Competition task baseline: Hi-VT5

- Textual representation
 - Hi-VT5 utilizes a spatial embedding to better align the layout information
- Visual representation
 - Hi-VT5 uses DiT to represent the page image as a set of patch embedding

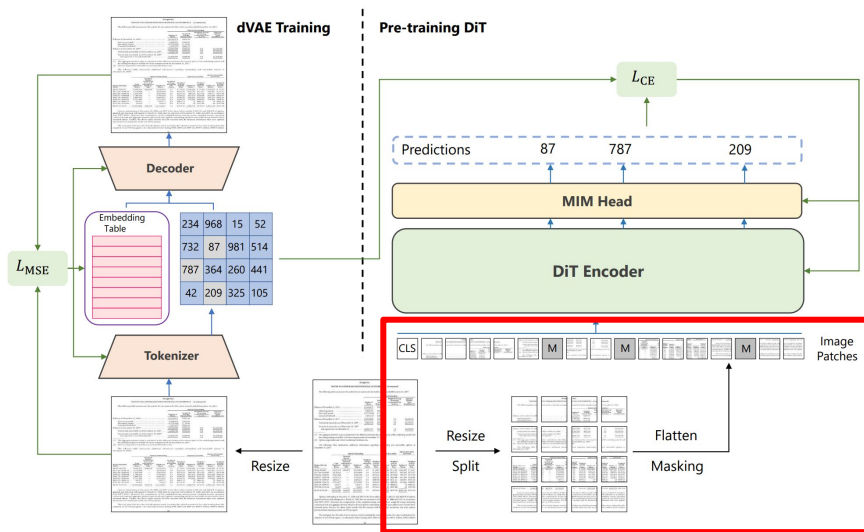


Figure 2: The model architecture of DiT with MIM pre-training.

Competition task baseline

Model	Size	Parameters	Max Seq. Length	Setup	Accuracy	ANLS	Ans. Page Accuracy
BERT [7]	Large	334M	512	Oracle	39.77	0.5904	100.00
				Max Conf.	34.78	0.5347	71.24
				Concat	27.41	0.4183	51.61
Longformer [1]	Base	148M	4096	Oracle	52.48	0.6177	100.00
				Max Conf.	45.87	0.5506	70.37
				Concat	43.91	0.5287	71.17
Big Bird [30]	Base	131M	4096	Oracle	55.31	0.6450	100.00
				Max Conf.	49.57	0.5854	72.27
				Concat	41.06	0.4929	67.54
LayoutLMv3 [9]	Base	125M	512	Oracle	58.81	0.6729	100.00
				Max Conf.	42.70	0.5513	74.02
				Concat	38.47	0.4538	51.94
T5 [20]	Base	223M	512	Oracle	59.00	0.6814	100.00
				Max Conf.	32.68	0.4028	46.05
				Concat	41.80	0.5050	—
Hi-VT5 (Ours)	Base	316M	20480	Oracle	50.01	0.6572	100.00
				Multipage	48.28	0.6201	79.23

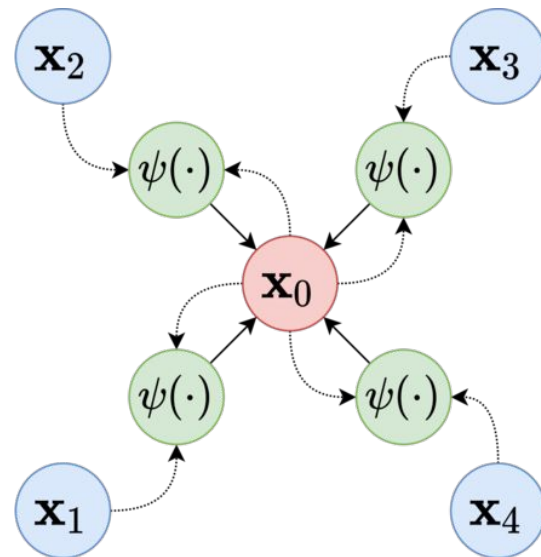
Methods

- Text-only KG
- Multi-modalities KG



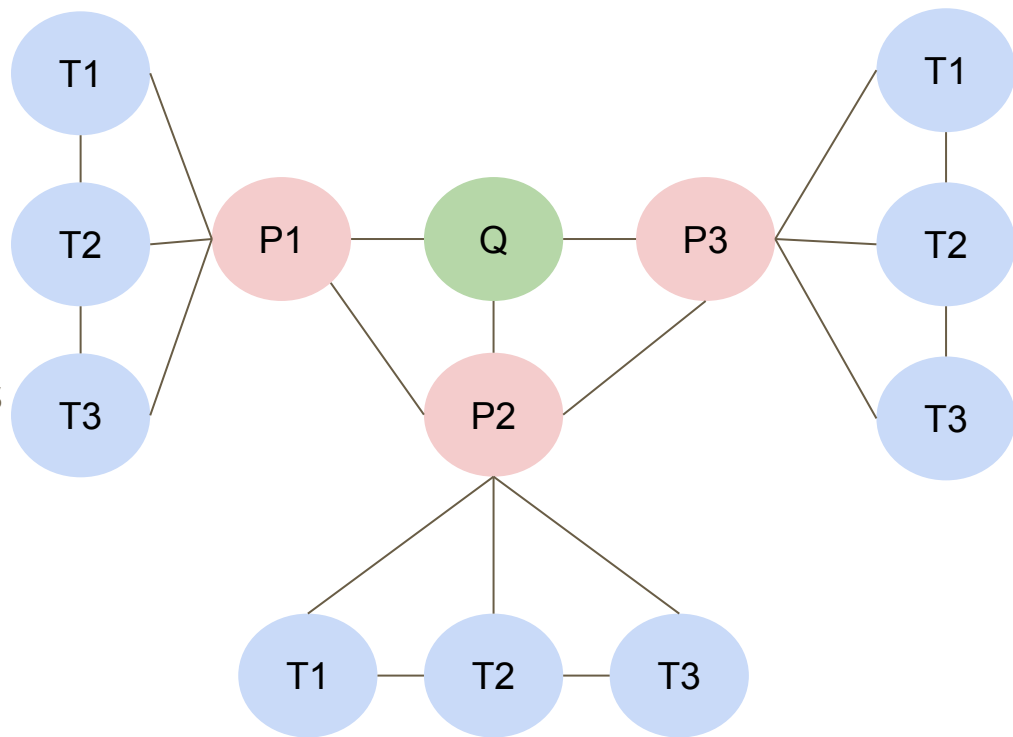
Method: Graph neural network (GNN)

- Graph convolutional network (GCN)
- Graph attention network (GAT)
 - GAT is just a different aggregation function with attention over features of neighbors, instead of a simple mean aggregation.



Method-1 Text-only KG

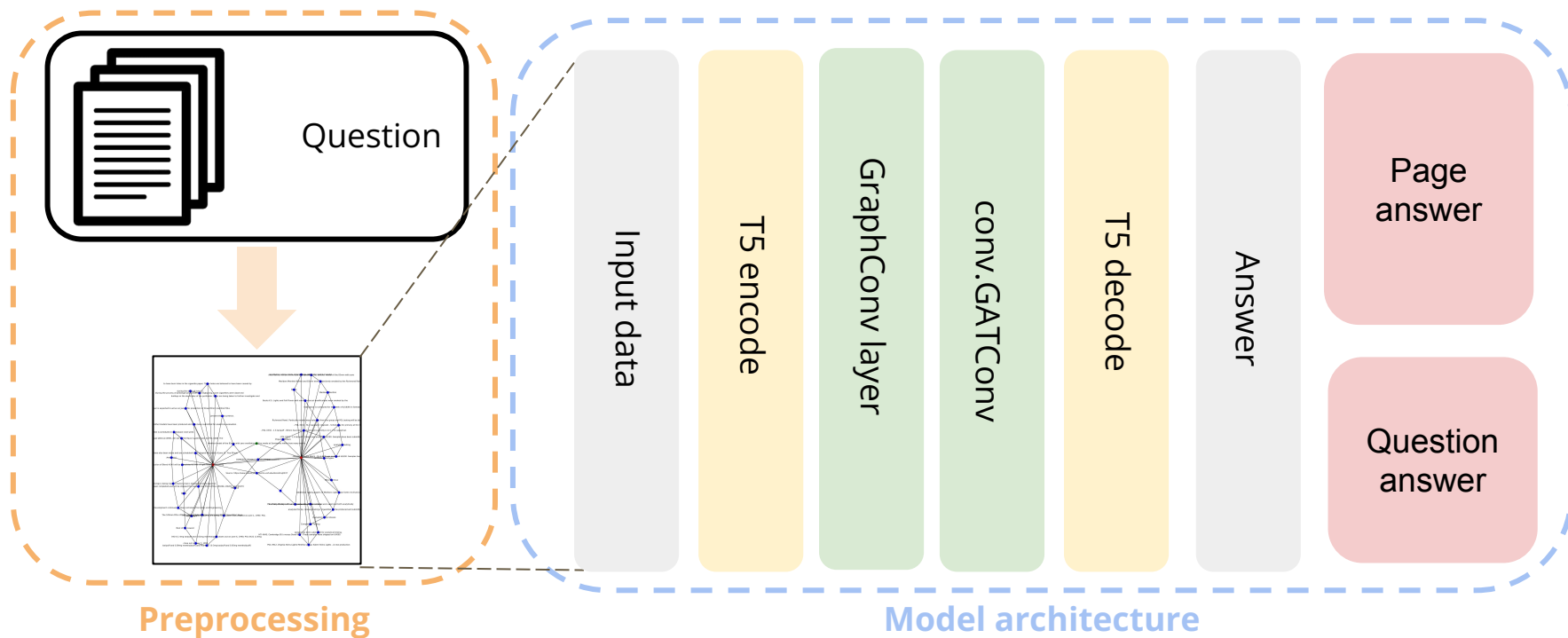
- KG elements:
 - Question node(Green)
 - Page nodes(Red)
 - Textual nodes(Blue)
- Edges:
 - Question node to Page nodes
 - Page node to page node
 - Page node to Textual nodes
 - Textual node to Textual node



- KG Construction

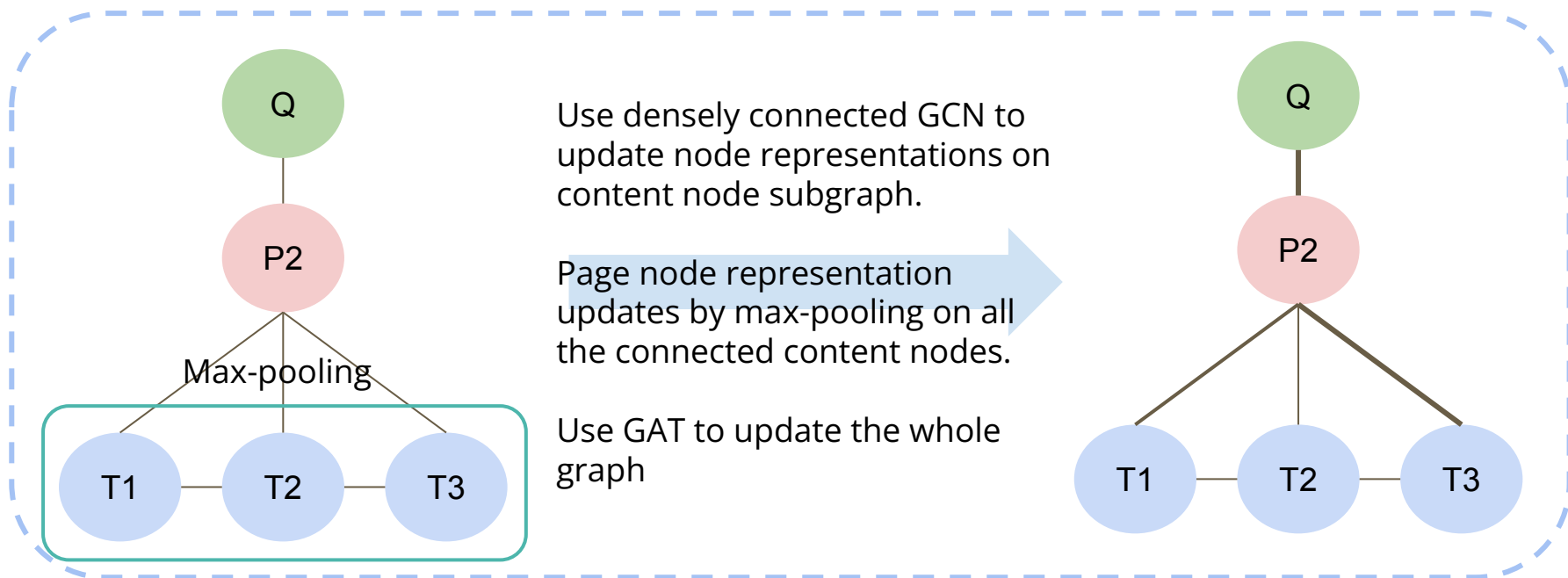
[illegible]

Method-1 Text-only KG



Method-1 Text-only KG

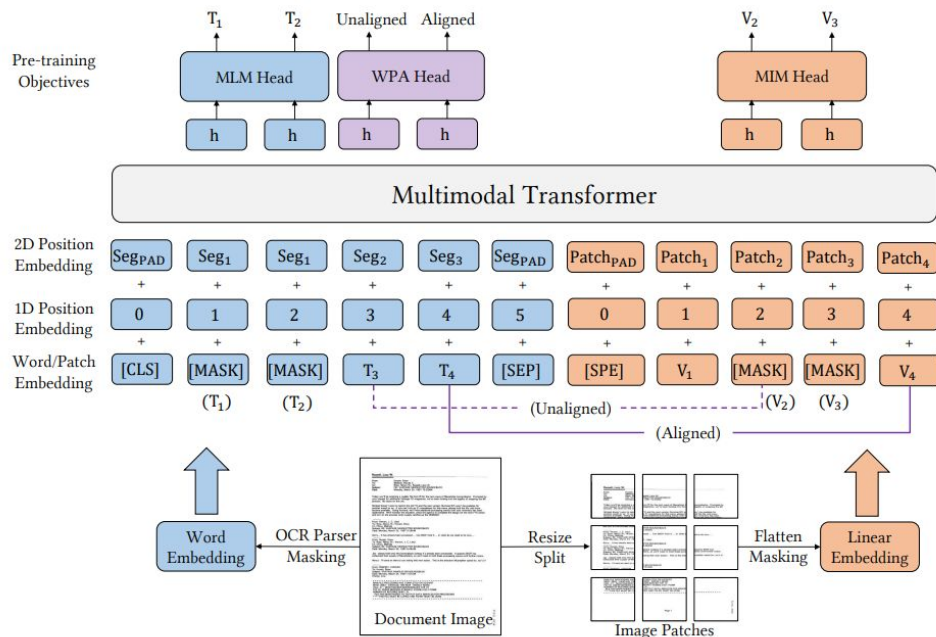
- Model architecture pipeline



Model architecture

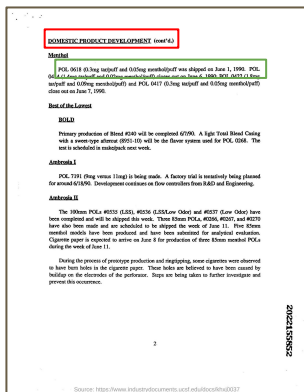
Method-2 Multi-modalities KG

- LayoutLMv3: pre-train multimodal Transformers for Document AI
- Text embedding:
 - A combination of word embeddings and position embeddings(OCR)
- Image embedding:
 - represent document images with linear projection features of image patches
 - Similar with DiT

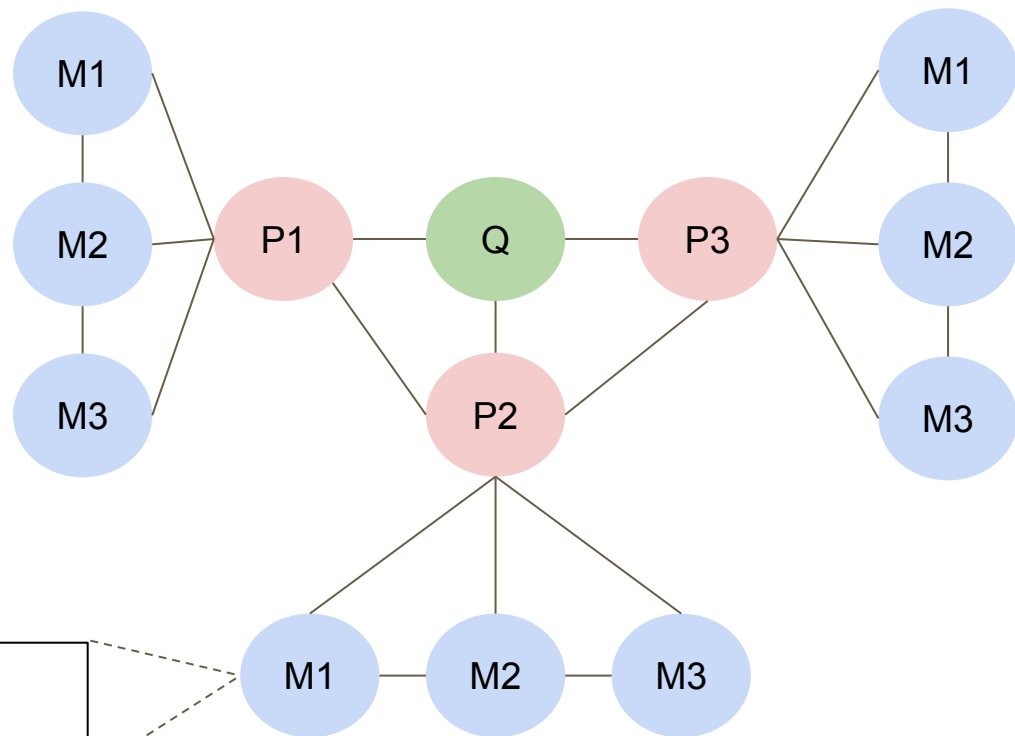


Method-2 Multi-modalities KG

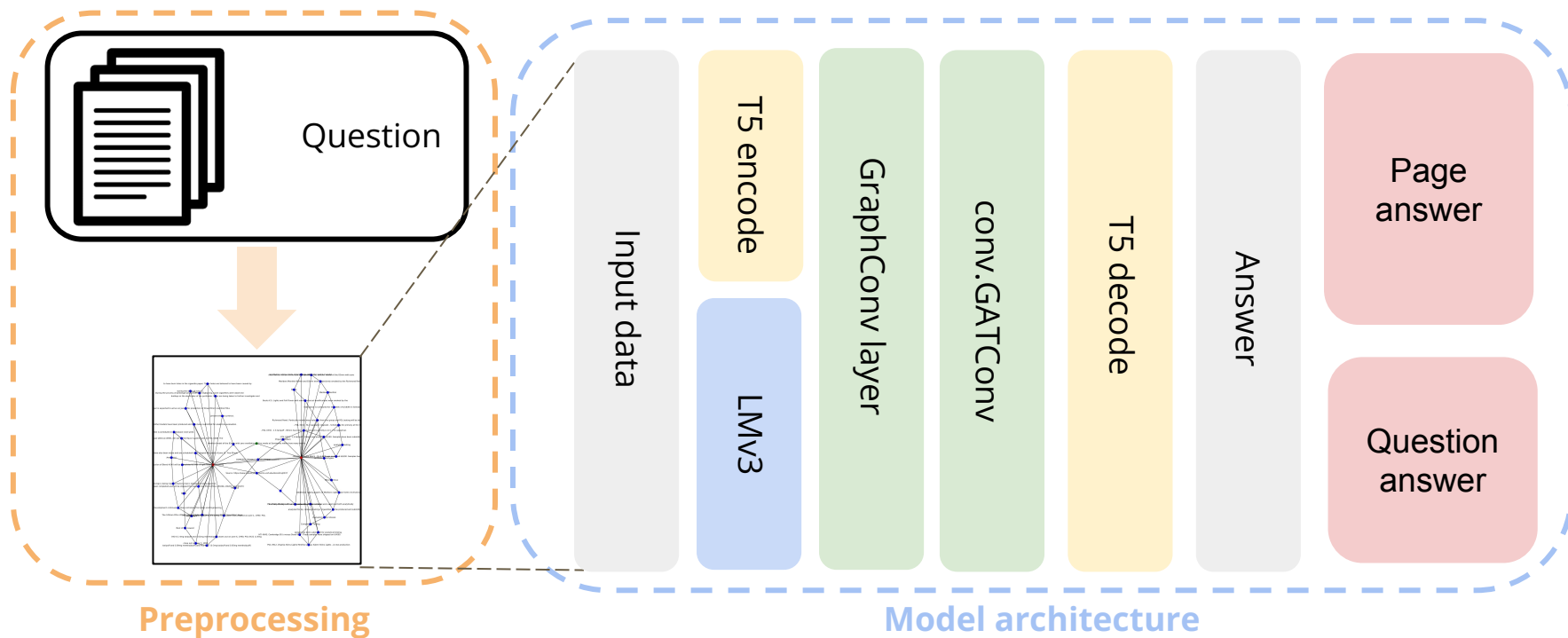
- KG nodes
 - Each node contains word, **bounding box**, **document image**
- KG edges
 - Same with method-1



1. Word text
2. OCR bbox
3. Doc image



Method-2 Multi-modalities KG



Results

- Still running
- We hope KG method is better than seq2seq method

Questions?