ale 2				OII	AL I	ILE		In	rior	• 1	141	DE	1 1111	101	LLO		J	110								1953
122	January		Yeb	CUATY.	Marc	ch	Apr	a	May	,	Jun		July	y	Augu	net:	Septes	mber	Octo	ber	Nove	mber	Decem	aber	Anne	ual
Station	Precipitation	Departure	Precipitation	Departure	Precipitation	Departure	Precipitation	Departure	Precipitation	Departure	Precipitation	Departure	Precipitation	Departure	Precipitation	Departure	Precipitation	Departure	Precipitation	Departure	Precipitation	Departure	Precipitation	Departure	Precipitation	Departure
DVANCE 5 655 CBANY LTON 6 ENE AITY NOERSON 1 5W	1.85	- 0.98 - 0.67 - 1.41 - 0.97	2.01		3.38 5.35 3.11	0.84	3.07 E 4.49 4.61	10000	3.15	1.70 0.16 1.71 3.38	2.45	- 2.18 - 2.84 - 2.71 - 5.46	.08- .43 1.74-	- 1.90 - 3.65 - 0.89 - 2.17	1.21-	2.89 2.39 3.65 2.36	1.30	- 2.43 - 4.25 - 1.75 - 2.65	1.05	- 0.36 - 1.75 - 1.57 - 0.85	1.42 E 2.06 1.61	- 2.12 - 0.29 - 1.07 - 2.43	1.69 E 1.63 2.55	STATE OF	21.30 (25.23 25.71	-12.70
MMAPOLIS ? SW PPLETON CITY SCADIA UTYASSE UN RANGER STA		0.83		7-0.80	4.18 3.90 4.52 4.05 6.03	0.75		- 2.22		2.10		- 1.10		2.24				- 1.45		- 0.15		1:70		200	25.00 26.13 222.37 27.12 24.15	-16.03
CLLEVIEW .	1.24	1.00	1.02		4.21		2.75	- 2.76	2.29-		3.16	- 1.44	-37	- 2.96	.92-	1.85	.00-	- 4.04	2.00	- 1.50	1.25	- 2.18	.61	- 2.20	19.83	-23.70

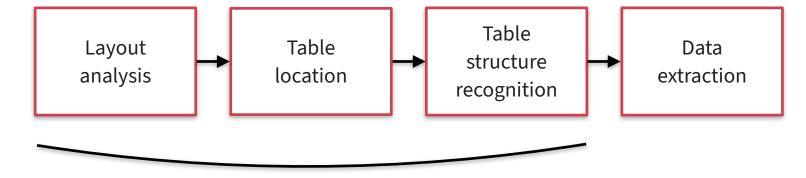
PDF Days Europe 2025

Extracting data from PDF tables

Fully explainable vs generative AI algorithms

Tamir Hassan, Co-Founder and CTO, Living PDF tamir@tamirhassan.com

Data extraction pipeline



Document understanding





Structure of the presentation

- Document understanding
 - Document image analysis
 - Object-level analysis of PDF
 - "Rule-based" systems
 - Machine learning
- Revisiting explainable algorithms
- Generative AI
- Comparison and conclusions







Document understanding



Document image analysis (1)

Separation of foreground from background



Thresholding or binarization

Source: Don Juan Archive, Vienna





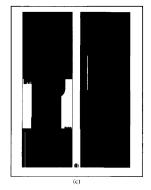


Document image analysis (2)

- Horizontal and vertical projection profiles
- Results in lines of text and blocks















Document image analysis (3)

Segmentation into individual characters often poses problems



Source: Ray Smith: An Overview of the Tesseract Engine, ICAR 2007





Document image analysis (4)

- Multiple segmentation hypotheses
- Segmentation-free and window-based recognition methods







Character classification

- A classifier, trained on ground-truthed data, returns the character code and corresponding confidence measure
- Typically convolutional neural networks are used
- Obtained words checked against the dictionary and/or linguistic rules
- Final result is the best combination of segmentation, classification and linguistic corrections





Processing PDF directly?

BT (begin text)

100 150 Td (move to coordinates)

[(T) 30 (his is a ty) 20 (pic) 40 (al sen) -10 (ten) -10 (ce.)] TJ

ET (end text)

[(AWAY again)]TJ AWAY again

[(A) 120 (W) 120 (A) 95 (Y again)]TJ AWAY again





Text fragments example

er etwas ist oder sein möchte Daimler-Benz, der achtet auf die Kleiderordnung; man trägt Blau im Schwabenkonzern, hell Fließband, dunkel auf der Führungs ebene. Und wer den Entscheidungsträ gern im Vorstand ganz nahe ist, der darf sich OFK-Mitglied nennen, der gehört zum oberen Führungskreis des Hauses

Ende Januar zogen rund 1000 der 1400 OFK Mitolieder in die Stuttgar-





Object-based layout analysis

UPDATE

In other news

- In former Soviet states, leaders watch uneasily as the call for democracy widens. Page 3
- Rebellion at La Scala ends as Riccardo Muti resigns as musical director after 19 years. Page 8
- A UN envoy says Syria has vowed to pull out all military and intelligence units from Lebanon by the end of April. Page 9
- Iraqi lawmakers elect a Sunni Arab as speaker of Parliament. Page 9

Daring to create

Mayor Bertrand Delanoë has worked hard to make his city a contemporary masterpiece. "Paris is a

museum, and that is a privilege," he says. "But if it wants to be loyal to its history, it needs to innovate, to dare — it needs to move into the 21st centu-

ry." Page 2



On the Web: www.iht.com





Das Windsor-Syndrom

Fordern Sie uns? Grillen Sie uns", appellierte Konzernehef Jürgen Schrempp an seine Führungskräfte. Auf dem Topmanagement-Meeting konterte er auch die Attacken seines Vorgüngers Edzord Reuter

Wer etwas ist oder sein michte bei Daireler-Honz, der achtot of the Kleiderordsung; man trapt. Der kantige Horst Zimmer (K), dem Sier "Chief Euconiver" (Japan, In-Blas im Schwabenkonzern, bell am Fliebhand, dunkel auf der Führungsbene. Und wer den Entscheidung strägern im Verstand gant nahe ist, der farf sich OPK-Mitglied nennen, der schört zum oberen Führungskreis des Hanses.

Ende Januar zogen rund 1000 der 1400 OPK-Mitglieder in die Stuttgaror Liederhalle nin, viele donkelblau gowandot and allo grepannt wie Chorlouber for einers großen Auftritt. Birgen Schrengp (53), der Chef. tatte gerufen, und er verfangte Manacement: "Dice ist unsere gemeinsame Veranstaltung, Nutzen Sie sie! Fordern Sie um ! Grillen Sie um !"

Für Spannung bei dem Treffen am 27. Januar war gosorge: Não zuvor lagen Markterfolg and Mismanagement to make beleimander; nie zuwor rankten tich so viele Gerüchte um Vorstände: nie zuvor hatte ein ehempliger Vorsitrender so mit dem Unternehmen abgeechnet wie jeiet Edward Reuter (70) n seinen Memoiren (siehe Kasten Spine 16%

Was sagt Schrenge inters our peinlichen Eleb-Panne, was zum veringlückten Smart? Welche Signali under der Vorstandschof in Richtung einer Kollegen Jürgen Hubbert (58 Pkw-Goschült) and Dieter Zetsche 44, friiher Entwicklung, heute Ver riob), die für das Debakel die Verantworking tragge.7

Der Manuschaftskapitän ließ sich sicht aus der Roserve lochere "Wir siton these kire als Team gegenüber", erkündete Schrenge. Ein größerei Revirement, so es denn dazu koment, bleibt bis nach der Hauptversummking

Democh gab as am Rande des Signistrations zwoi Toppersonalies: es als ornaipotentem Chef des Gochiftsbereichs Lige Europa iraner

larkt, and seine Division soll 1990 arn ersteamal wieder schwarze Zaben schreiben. Als Nachfelger wird der este Manu dos Barrichs Llow-Anrichostrang, Klaus Majer (44), go-

Durch rechtseitige Pensionierung stricht sich Peter Fietzek (59), Boeichsverstand und Asien-Beunftrager, den Folgen der Reorganisation in femost.

Künftig regionen in der Rugion dochina, Asean, China), die jeweils the governo Kongrerngeschäft von



iomlich egal was, wer unter ihm den Notefahregugvorstand abgab, verlijk den Konnern.

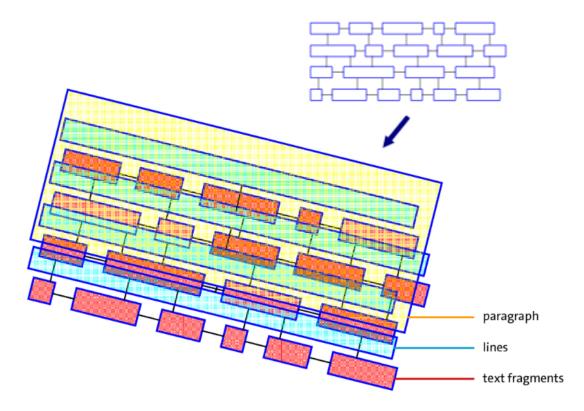
hen Gründen im Zenit seiner Kartiere: Die beiden neuen Modellreihen Actros (schwere Lastkraftwagen) ed Atogo (Verteller-Live) sind in

lagractures his Finanzdierstleistanon verant worken.

Auf der Veranstaltung selbst ging Zimmer geht aus gesandhoidi- is dann zunächst am einen Mann, der seit knapp drei Jahren zu internen-Daireler-Foren nicht mehr eingehalen wisch Edward Router, Day Wirken seies Vergingers kann Schrompp



Graph-based multi-level representation







Potential problems

- Text does not have to be in reading order
- Operators intended for kerning can be used to jump across columns
- Character encoding may be unknown
- Overprinting to simulate boldness

 and, of course, text may be part of an image object (or not recognized by OCR...)

Matrox product selection table

	Axio LE	Axio HD	Axio SD	RT.X2
Capture/editing formats				
<u> </u>				
HDV 1080i	K	Х	Х	Х
HDV 720p	k	X	У.	ngxt relaasa
DVCPRO HD	X.	χ	Х	
DV, DVCPRO, DVCAM	X	X	Х	Х
DVCPRO50	X	Х	Х	
P2 MXF DVCPRO50, DVCPRO HD	K	X	Х	
XDCAM MXF – DVCAM, IMX	X	X	У.	
XDCAM HD MXF, 18, 25, 35 mbps	X	X	Х	
Slow & Quick Motion				
MPEG-2 4:2:2 I-frame SD*	10-50 intops	10-50 mbps	10-50 mbps	10-25 mbps
Uncompressed 8-bit SD*	K	У.	X	





Table recognition (1/4)

Phosphorus (yellow or white)	7723-14-0	Only if it is a yellow or white form.
Sulfuric acid (acid acrosols including mists, vapors, gas, fog, and other airborne forms of any particle size)	7664-93-9	Only if it is an aerosol form as defined.
Vanadium (except when contained in an alloy)	7440-62-2	Except if it is contained in an alloy.
Zine (fume or dust)	7440-66-6	Only if it is in a fume or dust form.

The qualifier for the following three chemicals is based on the chemical activity rather than the form of the chemical. These chemicals are subject to EPCRA section 313 reporting requirements only when the indicated activity is performed.

Chemical/ Chemical Category	CAS Number	Qualifier
Dioxin and dioxin-like compounds (manufacturing; and the processing or otherwise use of dioxin and dioxin-like compounds if the dioxin and dioxin-like compounds are present as contaminants in a chemical and if they were created during the manufacture of that chemical.)	NA	Only if they are manufactured at the facility; or are processed or otherwise used when present as contaminants in a chemical but only if they were created during the manufacture of that chemical.
Isopropyl alcohol (only persons who manufacture by the strong acid process are subject, no supplier notification)	67-63-0	Only if it is being manufactured by the strong acid process. Facilities that process or otherwise use isopropyl alcohol are <u>not</u> covered.
Saccharin (only persons who manufacture are subject, no supplier notification)	81-07-2	Only if it is being manufactured.

There are no supplier notification requirements for isopropyl alcohol and saccharin since the processors and users of these chemicals are not required to report. Manufacturers of these chemicals do not need to notify their customers that these are reportable EPCRA





After candidate column finding (2/4)

Phosphorus (yellow or white)	7723-14-0	Only if it is a yellow or white form.
Sulfuric acid (acid acrosols including mists, vapors, gas, fog, and other airborne forms of any particle size)	7664-93-9	Only if it is an aerosol form as defined.
Vanadium (except when contained in an alloy)	7440-62-2	Except if it is contained in an alloy.
Zinc (fume or dust)	7440-66-6	Only if it is in a fume or dust form.

The qualifier for the following three chemicals is based on the chemical activity rather than the form of the chemical. These chemicals are subject to EPCRA section 313 reporting requirements only when the indicated activity is performed.

Chemical/ Chemical Category	CAS Number	Qualifier
Dioxin and dioxin-like compounds	NA	Only if they are manufactured at the
(manufacturing; and the processing or otherwise use		facility; or are processed or otherwise
of dioxin and dioxin-like compounds if the dioxin		used when present as contaminants in a
and dioxin-like compounds are present as		chemical but only if they were created
contaminants in a chemical and if they were created		during the manufacture of that chemical.
during the manufacture of that chemical.)		
Isopropyl alcohol (only persons who manufacture	67-63-0	Only if it is being manufactured by the
by the strong acid process are subject, no supplier		strong acid process. Facilities that process
notification)		or otherwise use isopropyl alcohol are not
•		covered
Caashawin (anh. nanana ayka manufaatina ana	81-07-2	Only if it is being manufactured
Saccharin (only persons who manufacture are	81-07-2	Only if it is being manufactured.
subject, no supplier notification)		

There are no supplier notification requirements for isopropyl alcohol and saccharin since the processors and users of these chemicals are not required to report. Manufacturers of these chemicals do not need to notify their customers that these are reportable EPCRA





After horizontal merging (3/4)

Phosphorus (yellow or white)	7723-14-0	Only if it is a yellow or white form.
Sulfuric acid (acid aerosols including mists, vapors, gas, fog, and other airborne forms of any particle size)	7664-93-9	Only if it is an aerosol form as defined.
Vanadium (except when contained in an alloy)	7440-62-2	Except if it is contained in an alloy.
Zinc (fume or dust)	7440-66-6	Only if it is in a fume or dust form.

The qualifier for the following three chemicals is based on the chemical activity rather than the form of the chemical. These chemicals are subject to EPCRA section 313 reporting requirements only when the indicated activity is performed.

Chemical/ Chemical Category	CAS Number	Oualifier
Chemical Chemical Category	CAS Number	<u>Onanner</u>
Dioxin and dioxin-like compounds	NA	Only if they are manufactured at the
(manufacturing; and the processing or otherwise use		facility; or are processed or otherwise
of dioxin and dioxin-like compounds if the dioxin		used when present as contaminants in a
and dioxin-like compounds are present as		chemical but only if they were created
contaminants in a chemical and if they were created		during the manufacture of that chemical.
during the manufacture of that chemical.)		
Isopropyl alcohol (only persons who manufacture by the strong acid process are subject, no supplier notification)	67-63-0	Only if it is being manufactured by the strong acid process. Facilities that process or otherwise use isopropyl alcohol are <u>not</u> covered.
Saccharin (only persons who manufacture are subject, no supplier notification)	81-07-2	Only if it is being manufactured.

There are no supplier notification requirements for isopropyl alcohol and saccharin since the processors and users of these chemicals are not required to report. Manufacturers of these chemicals do not need to notify their customers that these are reportable EPCRA





Final result (4/4)

all bottle forms of any particle size)		as defined.
Phosphorous (yellow or white)		Only if it is a yellow or white form.
Sulfuric acid (acid aerosols including mists, vapors, gas, fog, and other airborne forms of any particle size)	7664-93-9	Only if it is an aerosol form as defined.
Vanadium (except when contained in an alloy)	7440-62-2	Except if it is contained in an alloy.
Zinc (fume or dust)	7440-66-6	Only if it is in a fume or dust form.

The qualifier for the following three chemicals is based or the prical activity rather than the form of the chemical. These chemicals are subject to EPCRA section 313 reporting reporting

Chemical/ Chemical Category	ČAS Number	Qualifier
Dioxin and dioxin-like compounds (manufacturing; and the processing or otherwise use of dioxin and dioxin-like compounds if the dioxin and dioxin-like compounds are present as contaminants in a chemical and if they were created during the manufacture of that chemical.)	NA	Only if they are manufactured at the facility; or are processed or otherwise used when present as contaminants in a chemical but only if they were created during the manufacture of that chemical.
Isopropyl alcohol (only persons who manufacture by the strong acid process are subject, no supplier notification)	67-63-0	Only if it is being manufactured by the strong acid process. Facilities that process or otherwise use isopropyl alcohol are not covered.
Saccharin (only persons who manufacture are subject, no supplier notification)	81-07-2	Only if it is being manufactured.

There are no supplier notification requirements for isopropyl alcohol and saccharin since the processors and users of these chemicals are not required to report. Manufacturers of these chemicals do not need to notify their customers that these are





Limitations of conventional approaches

- "Knowledge gap:" Rules are inflexible and work in isolation on one level of granularity
- Difficult to backtrack

oft Case		G.F Holder IV (Hood V) *1			
_P1319	NC		(0)		
_P1319	NC		(0)		
_P1319	NC		(0)		
_P1214	NC		(0)		
_P1116	NC		(0)		
I P81 /	(0)		(U)		

ny and the ported. Congress ly six prise, but that hat the real Albert losed, subs of last

6 prisoner or suspecrovided by k after rereviewed vestigators tiries and criminal navy offistill under leaving court in New York on Tuesday after he was convicted of all charges.

technology The vero

CU	CURRENCIES New York						
	Tuesday 4 P.M.	Previous					
€1 =	\$1.3315	\$1.3369					
£1 =	\$1.9128	\$1.914					
\$1 =	¥104.49	¥104.925					
\$1 =	SF1.1645	SF1.1594					
Full	Full currency rates Page 14						
OII New York							







By James F

PARIS: T hind the cl London's v that who Deutsche E States wor biggest ma





Machine learning for structure detection

- Variety of approaches, mostly for images:
 - Object recognition (e.g. YOLO)
 - Pixel-based classification
 - Corner/feature detection
 - Pre-segmentation of image
 - Transformation of image
- Excel in simpler table structures and unclean documents
- Black box; difficult to patch; no direct link to PDF content stream
- Require huge dataset for training





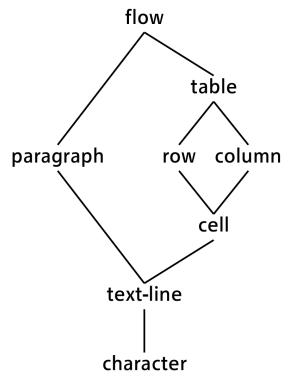


Revisiting explainable algorithms



Document analysis as an AI optimization problem

- Two components:
 - Document model with evaluation criteria (moderate)
 - Search through plausible interpretations (hard!)
- Which structures do we want to detect?
- Common to a wide variety of documents

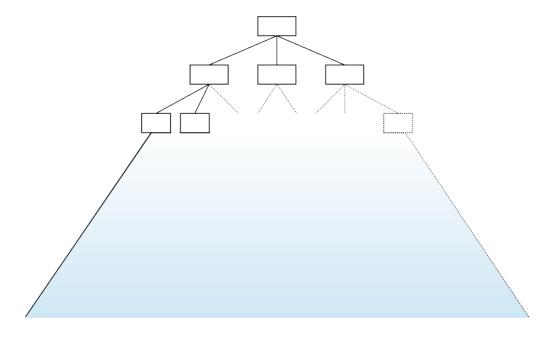






Guided search procedure (1/3)

Search space

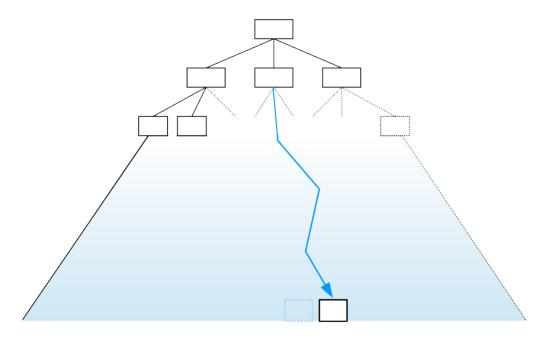






Guided search procedure (2/3)

Greedy search

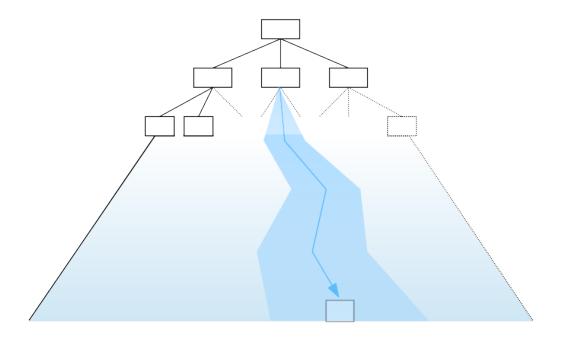






Guided search procedure (3/3)

Guided search







Rule-based vs AI optimization strategies

Rule-based

- Limited accuracy due to error propagation
- Easy to implement
- Fast
- Appropriate if source is known

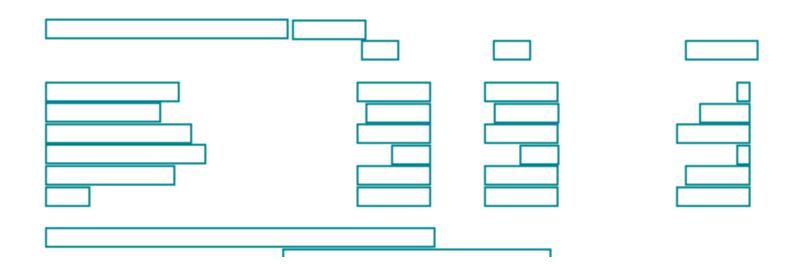
Al optimization

- Good accuracy for all documents supported by model
- Evaluation returns a confidence measure
- Implementation is complex
- Processing time is longer, depending on search strategy





Aside: No domain knowledge





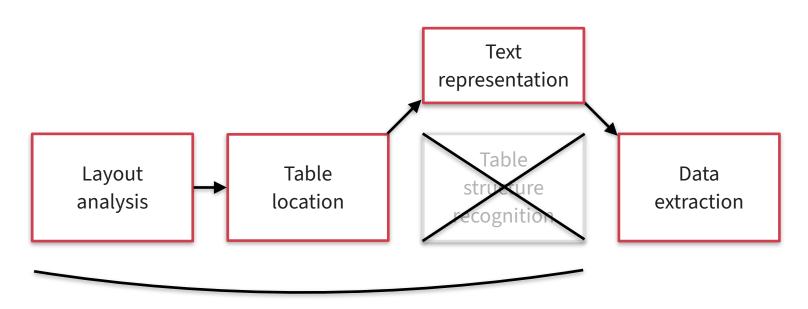




Generative Al



Generative AI option 1



Document understanding





Text representation

Differences with respect to Germany	Portugal	Greece	Spain	Italy	France
Q2 2006 – Q1 2010**	0.33	0.51	0.31	0.18	0.05
$\mathrm{Q1}\ 2009 - \mathrm{Q4}\ 2009$	-0.01	0.45	0.21	0.18	0.05
Q4 2009	0.17	0.70	0.26	0.09	-0.01
Q1 2010	0.64	0.72	0.56	0.43	0.25



Differences with respect to Germany	Portugal····	Greece	Spain	Italy···	France
Q2 2006 - Q1 2010**	0.33 · · · · · ·	0.51	0.31	0.18	0.05
Q1 · 2009 · - · Q4 · 2009 · · · · · · · · · · · · · · · · · ·	-0.01	0.45	0.21	0.18	0.05
Q4·2009······	0.17	0.70	0.26	0.09	-0.01
Q1 · 2010 · · · · · · · · · · · · · · · · · ·	0.64	0.72	0.56	0.43	0.25

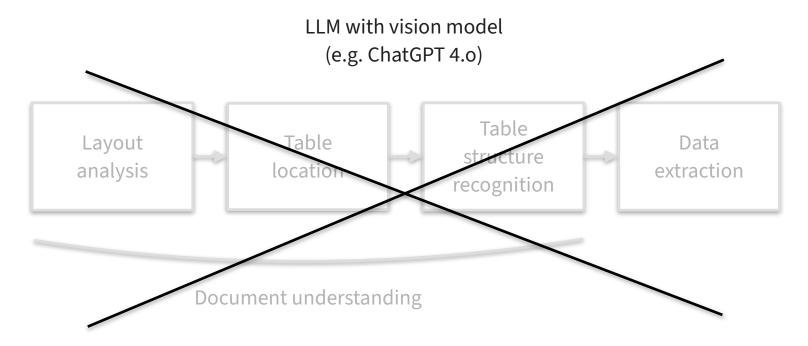
"According to the following table, what was the figure for Spain in Q4 2009?"

Source: ICDAR 2013 Table Competition Dataset





Generative AI option 2







Vision model

Differences with respect to Germany	Portugal	Greece	Spain	Italy	France
Q2 2006 – Q1 2010**	0.33	0.51	0.31	0.18	0.05
$\mathrm{Q1}\ 2009 - \mathrm{Q4}\ 2009$	-0.01	0.45	0.21	0.18	0.05
Q4 2009	0.17	0.70	0.26	0.09	-0.01
Q1 2010	0.64	0.72	0.56	0.43	0.25



"According to the following table, what was the figure for Spain in Q4 2009?"

"Please output the attached table in HTML format."

Here is the HTML representation of the table shown in your image:

```
html
<thead>
 Differences with respect to Germany
 Portugal
 Greece
 Spain
 Italy
 France
 </thead>
Q2 2006 - Q1 2010**
 0.33
 0.51
 0.31
 0.18
 0.05
                        \downarrow
```





Explainable vs generative AI strategies

Al optimization (explainable)

- Good accuracy for all documents supported by model
- Fully explainable result
- Evaluation returns a confidence measure
- Implementation is complex
- Typically requires less processing power than GenAl
- Typically no domain knowledge

Generative AI (RAG with vision model)

- Generally good accuracy for all documents typically found on the Web
- Generally no confidence measures
- Easy to implement
- May require greater processing power than AI search
- May hallucinate; can be combined with other approaches to detect such cases





Combining GenAI with explainable methods

- Resort to generative AI when explainable methods fail to find a suitable result (with sufficient confidence)
 - might be stuck in a local maximum
- Use the explainable method's evaluation component to verify the plausibility of the GenAI result







Thank you

Tamir Hassan tamir@tamirhassan.com

