



*Integral* > Business  
Intelligence

AI for business. In-House. Private. Secure.

# About Integral Business Intelligence



Tony Sclafani  
Co-Founder / Partner

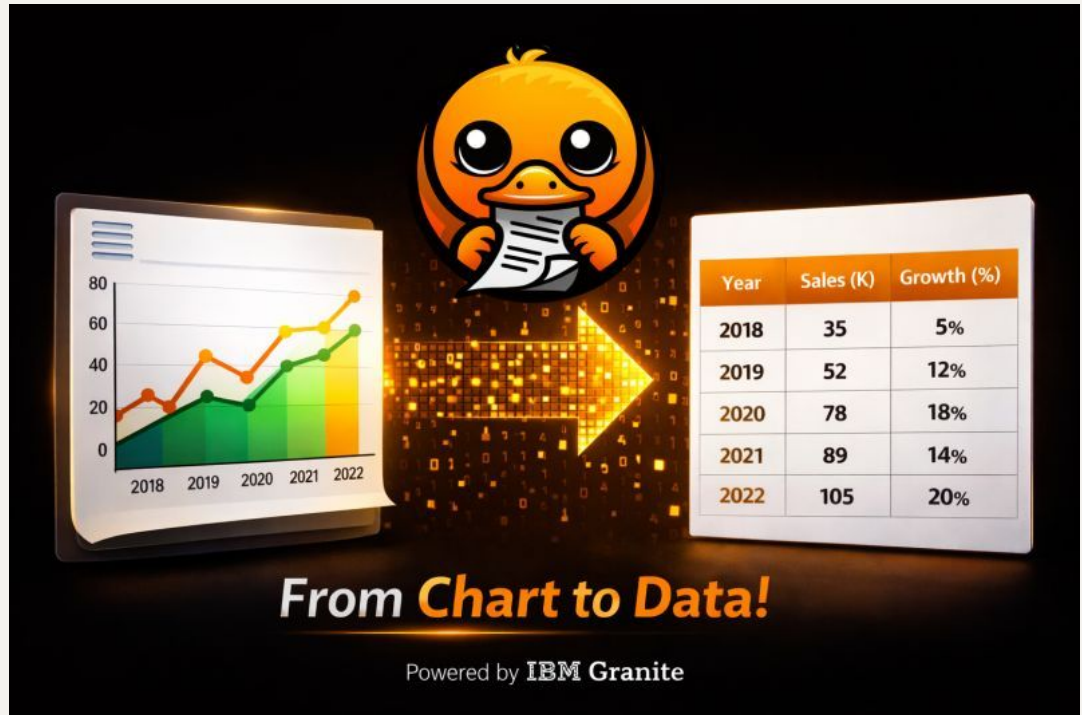


<https://integralbi.ai>

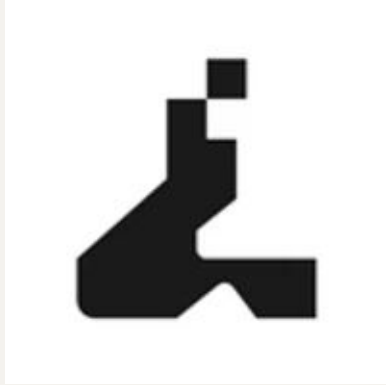
- ✓ AI development & systems integration
- ✓ Ideal customers: small and medium businesses with sensitive data
- ✓ Extract, transform, load (ETL) sensitive data into in-house AI platforms
  - ✓ Paper-to-Digital
  - ✓ Digital-to-Database

# In-House AI for Client-Confidential Data

- Text
- Document structure
- Tables
- Charts
- Images



# In-House AI for Client-Confidential Data



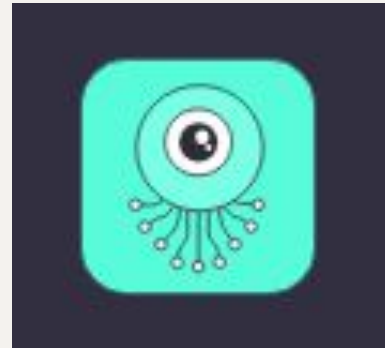
Marker  
by datalab.to



Docling  
by IBM

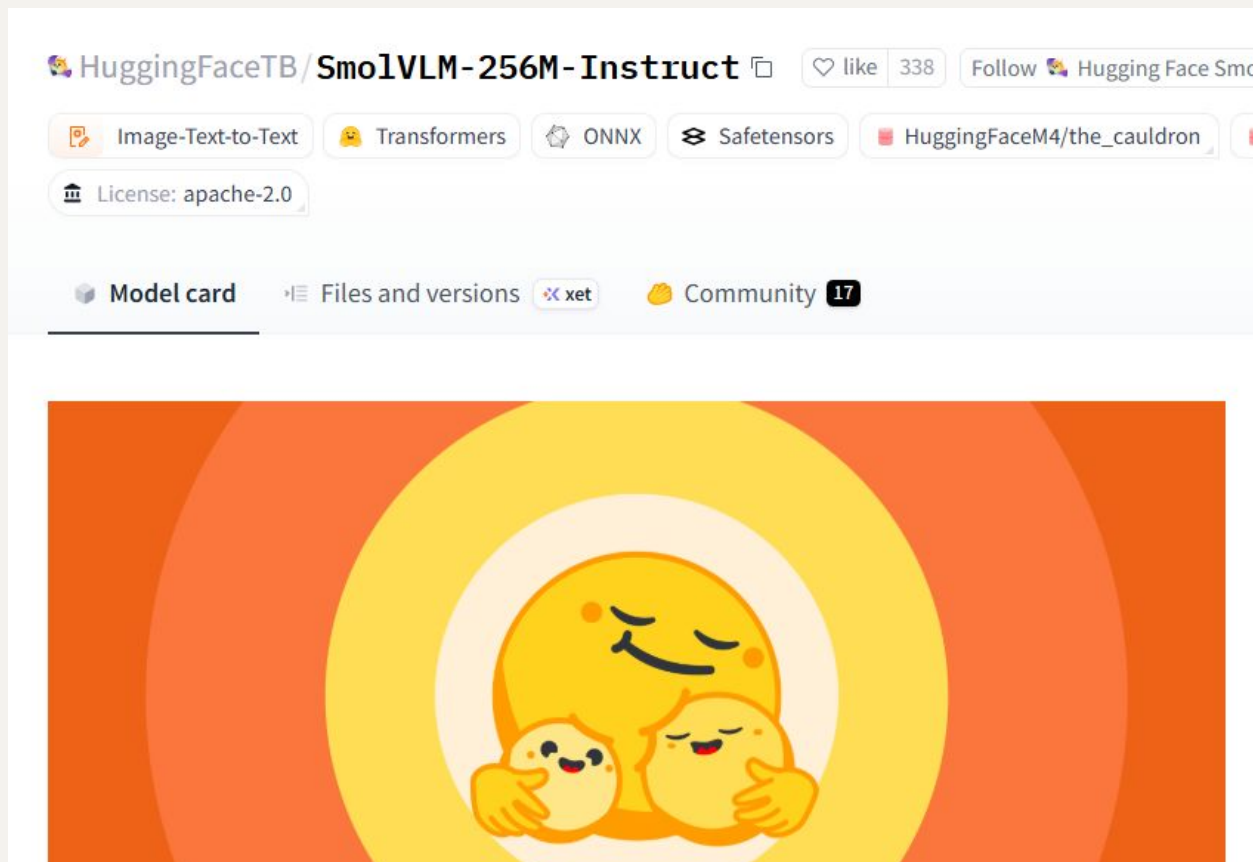


DeepSeek-OCR  
by DeepSeek

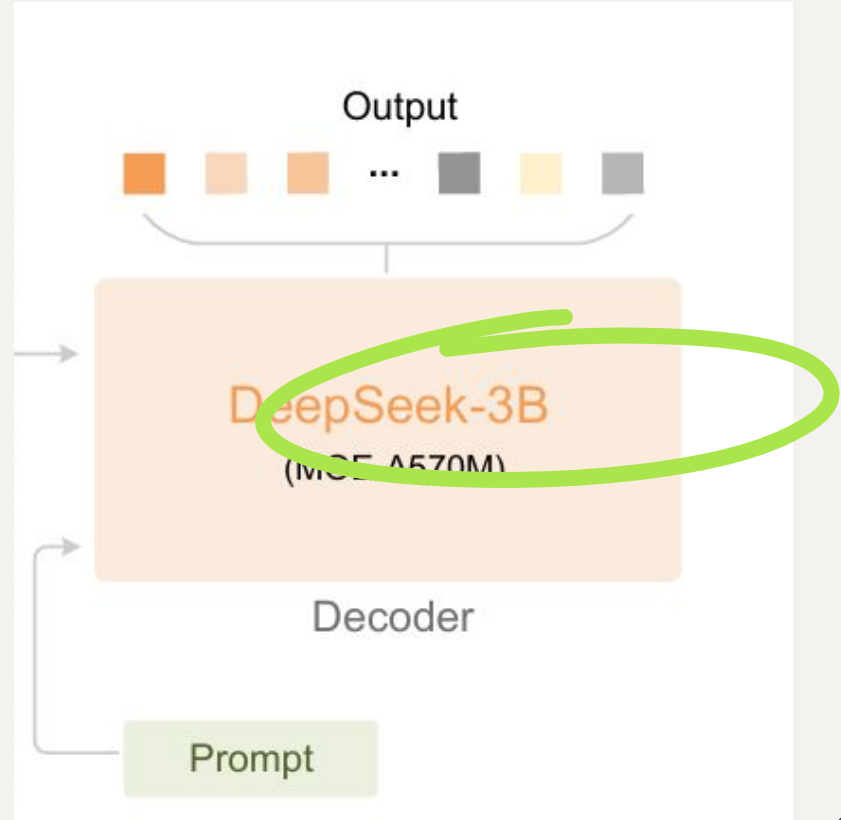


Kreuzberg

# Limitations



# Overcoming Limitations



# New Developments



Model	Release	Research Purpose
DeepSeek-OCR	October 21, 2025	Compression
DeepSeek-OCR-2	January 28, 2026	Reading Order



## DeepSeek-OCR: Contexts Optical Compression

Haoran Wei, Yaofeng Sun, Yukun Li

DeepSeek-AI

### Abstract

We present DeepSeek-OCR as an initial investigation into the feasibility of compressing long contexts via optical 2D mapping. DeepSeek-OCR consists of two components: DeepEncoder and DeepSeek3B-MoE-A570M as the decoder. Specifically, DeepEncoder serves as the core

“Our work represents an initial exploration into the boundaries of vision-text compression, investigating how many vision tokens are required to decode  $N$  text tokens.”

800 vision tokens. In production, DeepSeek-OCR can generate training data for LLMs/VLMs at a scale of 200k+ pages per day (a single A100-40G). Codes and model weights are publicly accessible at <http://github.com/deepseek-ai/DeepSeek-OCR>.



# DeepSeek-OCR

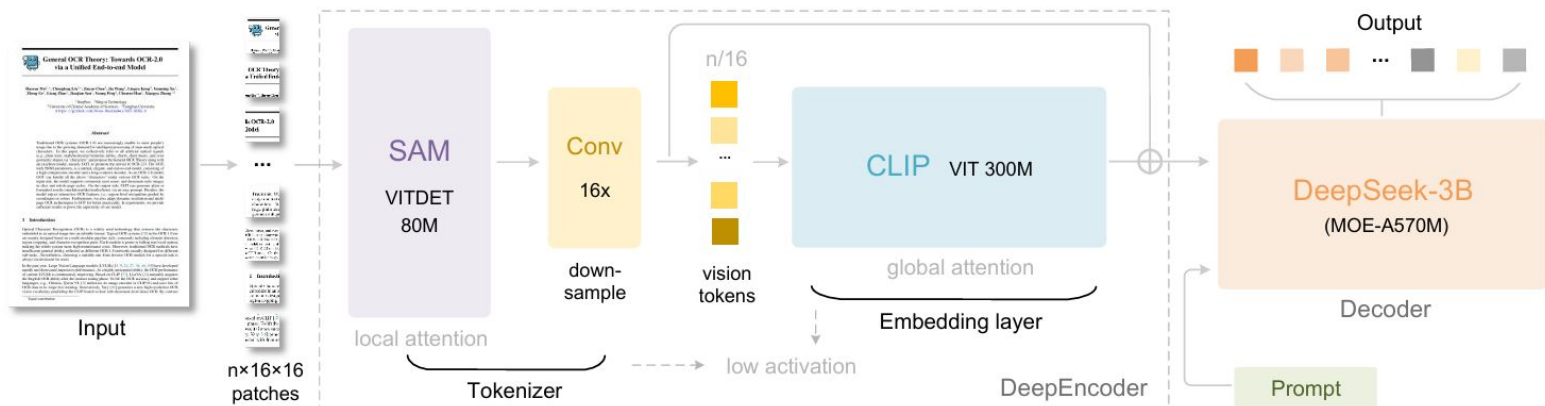
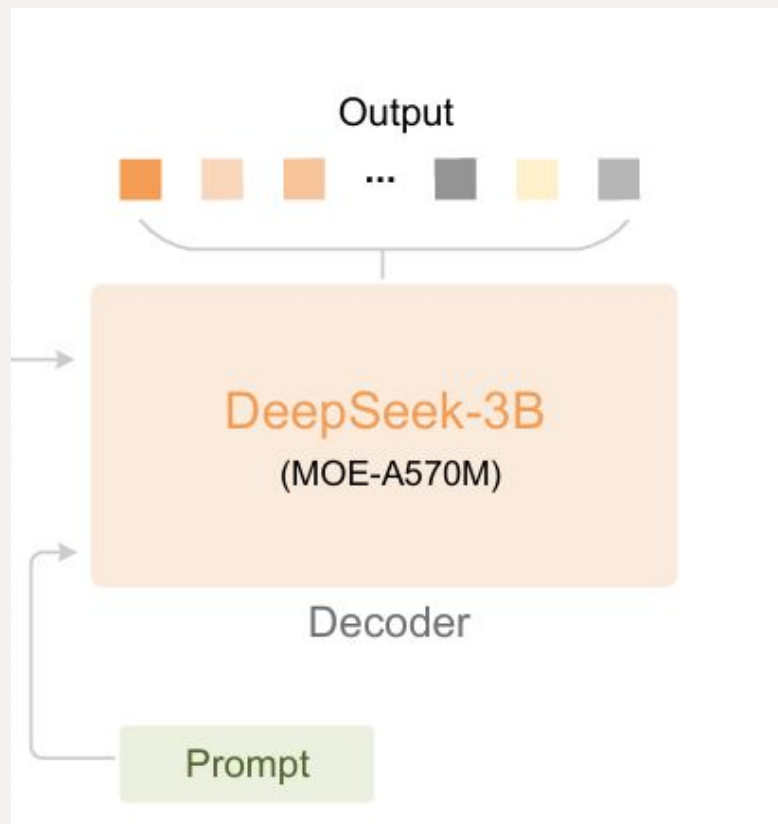


Figure 3 | The architecture of DeepSeek-OCR. DeepSeek-OCR consists of a DeepEncoder and a DeepSeek-3B-MoE decoder. DeepEncoder is the core of DeepSeek-OCR, comprising three components: a SAM [17] for perception dominated by window attention, a CLIP [29] for knowledge with dense global attention, and a  $16 \times$  token compressor that bridges between them.

## 2 Stage Training - DeepSeek-OCR

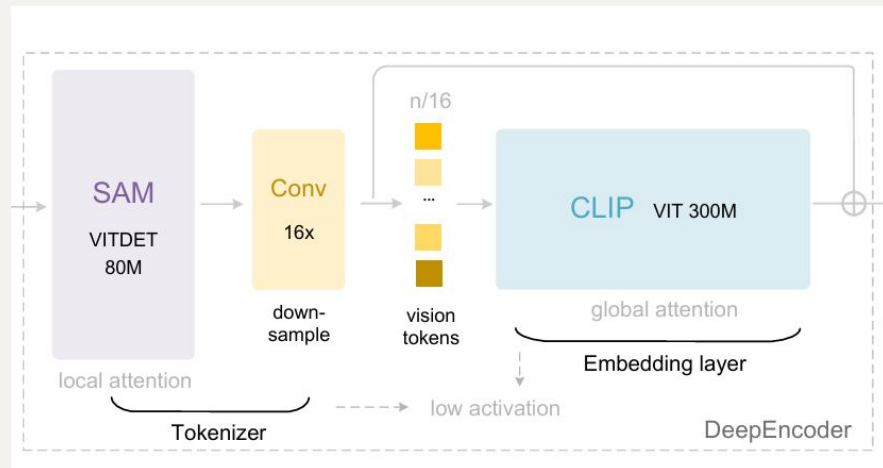
- Stage 1 - DeepEncoder
- Stage 2 - “Full Model” / Full Pipeline

# Decoder



# DeepEncoder

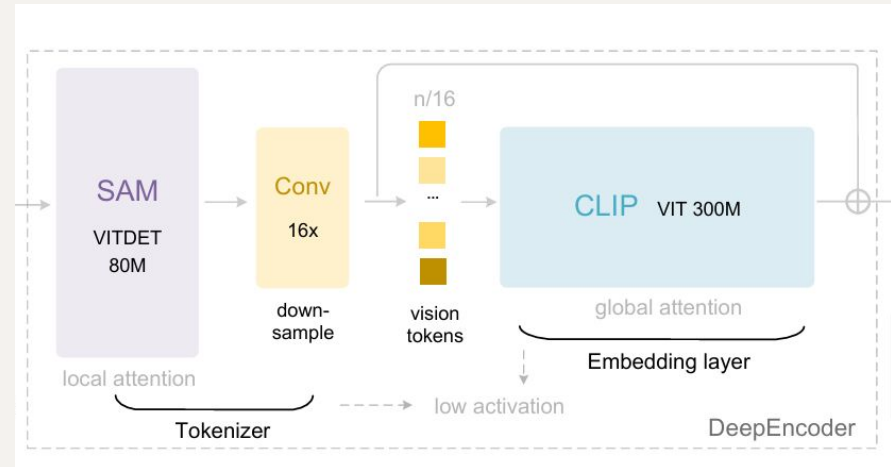
- SAM
  - Perception ("where are things?")
  - Image-to-embeddings
  - 80M params
  - Local attention
  - 1024×1024 image yields 4096 patch tokens
- 16× Convolutional Compressor
  - Token reduction
  - Hard-coded 16x
  - Embeddings-to-less-embeddings
- CLIP
  - Knowledge enrichment ("what do things mean?")
  - Global attention
  - Embeddings-to-embeddings



# Compression

- SAM
  - 1024×1024 image
  - 16px tiles
  - 64 patches x 64 patches
  - 1 token per patch
  - Total 4096 patch tokens
- 16× Convolutional Compressor
  - Token reduction
  - Hard-coded 16x
  - Compress 4096 tokens to 256

1. Are the output tokens any good?
2. Input resolution is user-adjustable



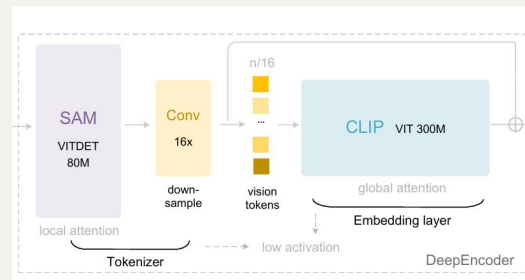
# Are the output tokens any good?

Table 2 | We test DeepSeek-OCR's vision-text compression ratio using all English documents with 600-1300 tokens from the Fox [21] benchmarks. Text tokens represent the number of tokens after tokenizing the ground truth text using DeepSeek-OCR's tokenizer. Vision Tokens=64 or 100 respectively represent the number of vision tokens output by DeepEncoder after resizing input images to 512×512 and 640×640.

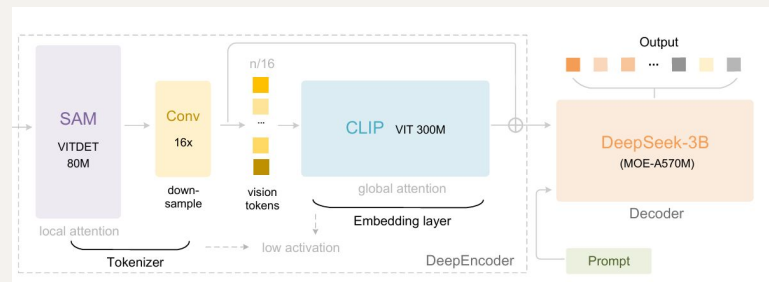
Text Tokens	Vision Tokens =64		Vision Tokens=100		Pages
	Precision	Compression	Precision	Compression	
600-700	96.5%	10.5×	98.5%	6.7×	7
700-800	93.8%	11.8×	97.3%	7.5×	28
800-900	83.8%	13.2×	96.8%	8.5×	28
900-1000	85.9%	15.1×	96.8%	9.7×	14
1000-1100	79.3%	16.5×	91.5%	10.6×	11
1100-1200	76.4%	17.7×	89.8%	11.3×	8
1200-1300	59.1%	19.7×	87.1%	12.6×	4

## 2 Stage Training - DeepSeek-OCR

- Stage 1 - DeepEncoder
  - ☒ SAM: Training
  - ☒ Compressor: Training
  - ☒ CLIP: Training
  - ☒ DeepSeek-3B-MoE decoder: Omitted



- Stage 2 - "Full Model" / Full Pipeline
  - ☒ SAM: Frozen
  - ☒ Compressor: Frozen
  - ☒ CLIP: Training
  - ☒ DeepSeek-3B-MoE decoder: Training



# Research Conclusions

- ...with a 10× compression ratio, the model's decoding precision can reach approximately 97%
- When compressing tokens by nearly 20×, we find that precision can still approach 60%



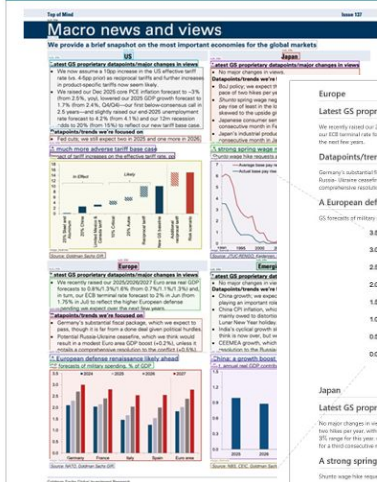
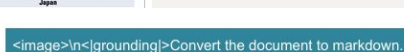
## Applications - DeepSeek-OCR

- OCR
- LLM/VLM pretraining
- “Deep Parsing”

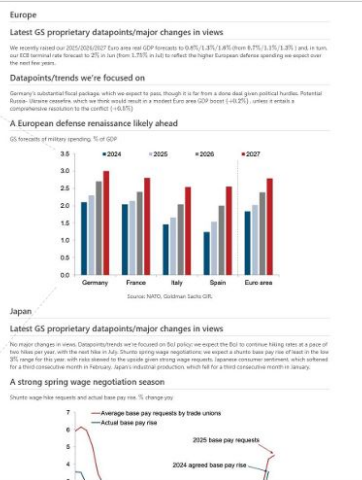
# Deep Parsing



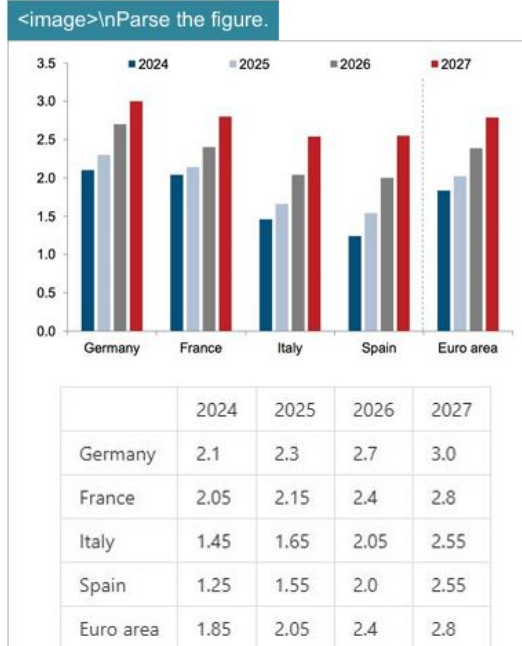
Input image



## Result



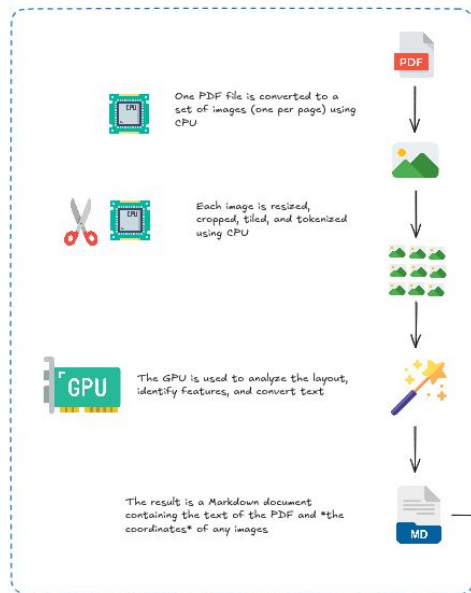
## Rendering



## Deep Parsing

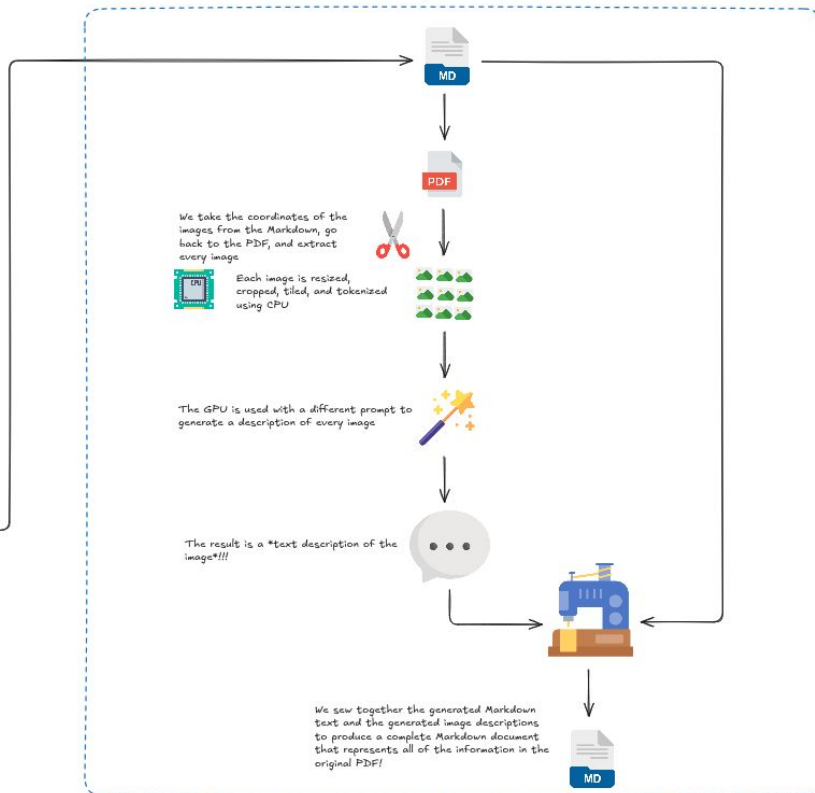
# Deep Parsing

## DeepSeek Single File Pipeline Pass 1



Looping this pipeline is inefficient. That would only load the images from one file at a time into the GPU - poorly utilizing the compute and slowing the process.

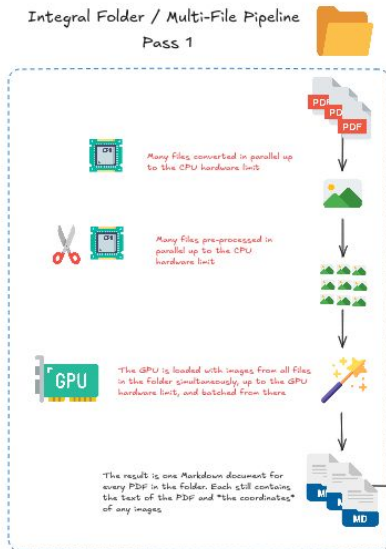
## Pass 2 (Implied by DeepSeek) "Deep Parsing"



# Deep Parsing

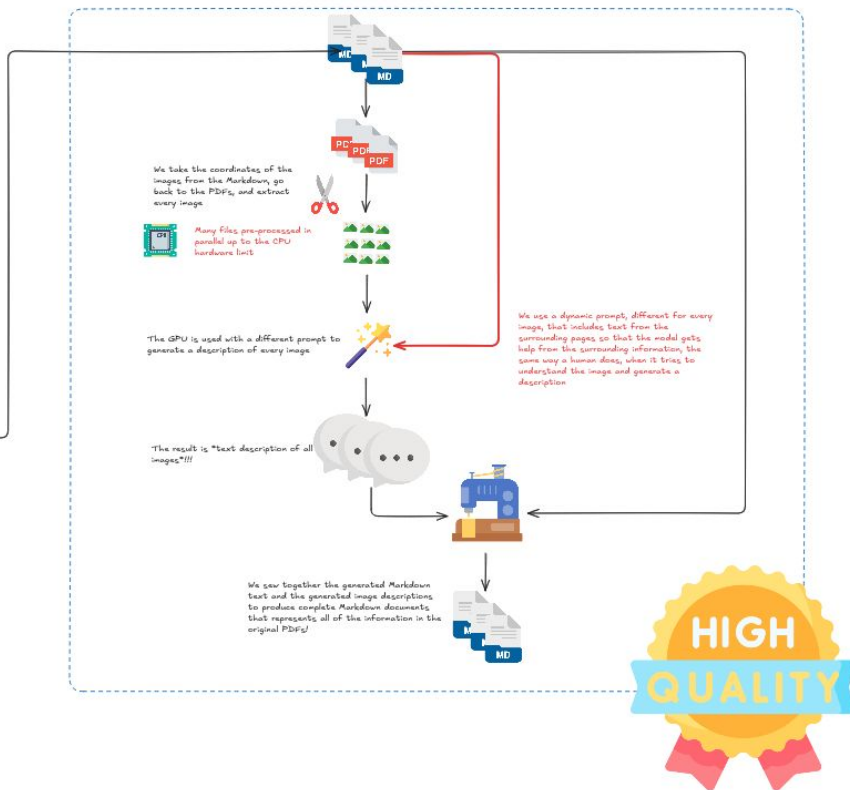
Integral Folder / Multi-File Pipeline

Pass 1



Integral Folder / Multi-File "Deep Parsing" Pipeline

CONTEXT AWARE Pass 2



# DeepSeek-OCR Companion Docs

<https://github.com/integral-business-intelligence/deepseek-ocr-companion>

A screenshot of a YouTube video player. The video title is "Performance on a Page". The video content displays two columns of performance metrics. The left column, titled "Performance", lists: "7 to 9 seconds/page", "30+ hours per gigabyte of documents", "vLLM 20x faster than HF Transformers", "7 GB model weights and 83 GB VRAM usage", and "Hallucinations were observed". The right column, titled "Real-World Setup", lists: "US patent application PDFs (text + image)", "Requires two-pass 'Deep Parsing'", "Single NVIDIA H100 GPU", "vLLM", and "3 resolution tests: Tiny, Base, Gundam". The video player interface includes the YouTube logo, a search bar, and a channel name "Integral Business Intelligence" with 7 subscribers. The video title "DeepSeek-OCR Speed Test" is visible at the bottom left. The DeepSeek logo is in the bottom right corner of the video frame. Interaction buttons for likes, shares, and saves are at the bottom right of the player.

YouTube

Search

### Performance on a Page

Performance	Real-World Setup
<ul style="list-style-type: none"><li>7 to 9 seconds/page</li><li>30+ hours per gigabyte of documents</li><li>vLLM 20x faster than HF Transformers</li><li>7 GB model weights and 83 GB VRAM usage</li><li>Hallucinations were observed</li></ul>	<ul style="list-style-type: none"><li>US patent application PDFs (text + image)</li><li>Requires two-pass "Deep Parsing"</li><li>Single NVIDIA H100 GPU</li><li>vLLM</li><li>3 resolution tests:<ul style="list-style-type: none"><li>Tiny</li><li>Base</li><li>Gundam</li></ul></li></ul>

deepseek

DeepSeek-OCR Speed Test

Integral Business Intelligence  
7 subscribers

Subscribe

1 | Share | Save



## DeepSeek-OCR 2: Visual Causal Flow

Haoran Wei, Yaofeng Sun, Yukun Li

DeepSeek-AI

### Abstract

Inspired by this cognitive mechanism, DeepEncoder V2 is designed to endow the encoder with causal reasoning capabilities, enabling it to intelligently reorder visual tokens prior to LLM-based content interpretation.

through two-cascaded 1D causal reasoning structures, thereby offering a new architectural approach with the potential to achieve genuine 2D reasoning. Codes and model weights are publicly accessible at <http://github.com/deepseek-ai/DeepSeek-OCR-2>.

# New Developments



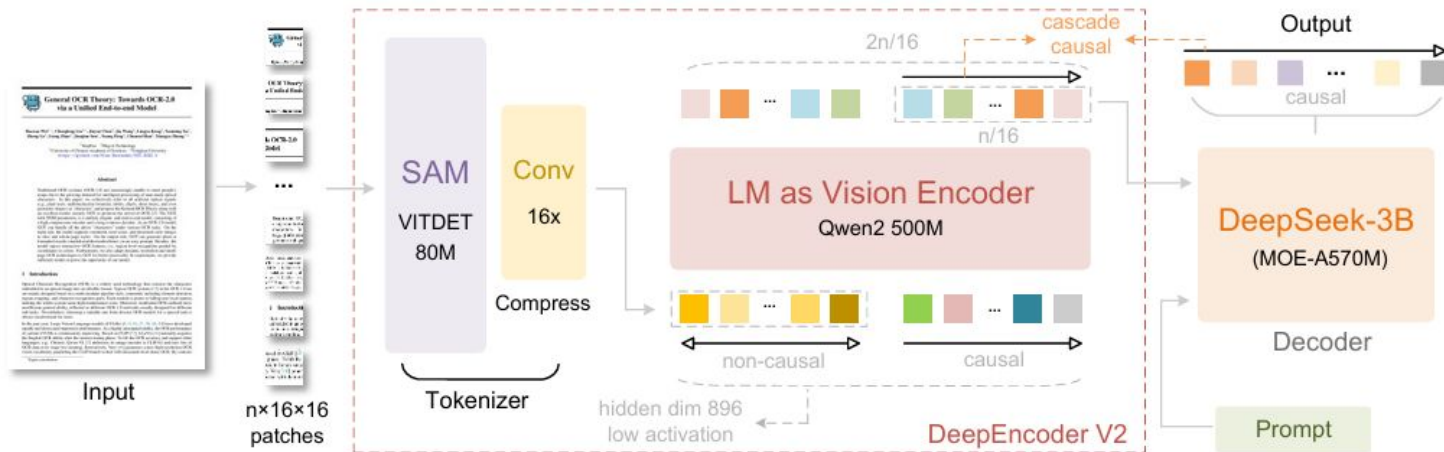
Model	Release	Research Purpose
DeepSeek-OCR	October 21, 2025	Compression
DeepSeek-OCR-2	January 28, 2026	Reading Order

# Research Hypothesis

- The human visual system ... functions as visual tokens, locally sharp yet globally aware.
- ...unlike existing encoders that rigidly scan tokens from top-left to bottom-right...
- ...visual tokens in models should be selectively processed with ordering highly contingent on visual semantics rather than spatial coordinates
- ... fundamentally reconsider the architectural design of vision language models (VLMs), ... Directly flattening image patches in a predefined raster-scan order introduces unwarranted inductive bias that ignores semantic relationships.



# DeepSeek-OCR-2

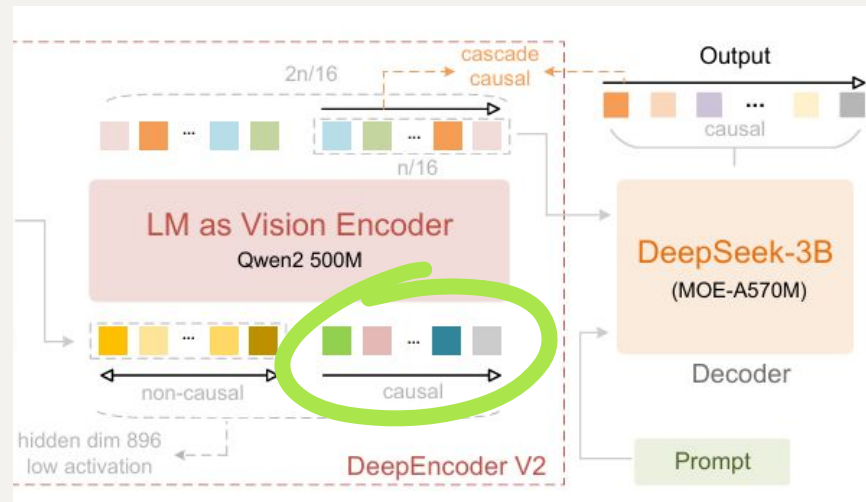


Replaced CLIP with Qwen2-0.5B.  
Reused 3B decoder model.

token compression mechanism from DeepEncoder, a vision encoder that reduces visual tokens by a factor of 16. By replacing the encoder's CLIP module with a compact language model on masks, this LM-style vision encoder acquires the ability to initiate causal modeling of visual sequences.

# 3 Stage Training - DeepSeek-OCR-2

- Stage 1 - DeepEncoder v2
  - ☒ SAM: Training
  - ☒ Compressor: Training
  - ☒ Qwen Model: Training
  - ☒ DeepSeek-3B-MoE decoder: Omitted
- Stage 2 - Query Enhancement / Full Pipeline
  - ☒ SAM: Frozen
  - ☒ Compressor: Frozen
  - ☒ Qwen Model: Training
  - ☒ DeepSeek-3B-MoE decoder: Training
- Stage 3 - "Full Model" / Full Pipeline
  - ☒ SAM: Frozen
  - ☒ Compressor: Frozen
  - ☒ Qwen Model: Frozen
  - ☒ DeepSeek-3B-MoE decoder: Training



Called "learnable queries" which can attend to all non-causal/vision tokens + previous queries + itself

Stage 2 learns reading order

# DeepSeek-OCR-2 Research Conclusions

- Replacing CLIP with an LLM works
- OmniDocBench score of 91% versus 87% for DeepSeek-OCR

# Research Conclusions

- DeepSeek-OCR serves two primary production use cases:
  - An online OCR service that reads image/documents for DeepSeek-LLMs
  - A pretraining data pipeline that performs batch PDF processing

# *Integral* > Business Intelligence



Tony Sclafani  
Co-Founder / Partner



<https://integralbi.ai>

