

	Sid	Japanese POC Differential parsing in MLLM for rendering and parsing (Binary + Ocr)		Branch	Main	
				Yash	feature_yash_attack1	feature_name_focus
				PB	feature_pb_injection3	
				Sid	feature_sid_textonimage poc	
				Ash	feature_ash_setup	
		Framework	env			
		LangGraph	uv/conda			
Input	PDF	[Finance, Education ,Health]				
		60	75	94	90	
Step 1	PDF Processing and extraction	[PymuPDF/ReportLab(Direct Binary Parsing)]	[Tesseract/Docling(OCR )  [Ulamaparser/Landing AI/deepseek/paddleocr(Agentic Document Extraction	[Gemini, Claude, Gpt] (MLLM)  	Function/Tools calls	
		C1 - Content extracted	C2	C3	C4	
Step 2	Analyze					
		Summary, Domain, Intended Task, Sub Tasks, Evidence for task, Pre-conditions, effects, Field information, Conatins -> images, tables, etc., B-box from step 1), Original Document Source (PPT, ARXIV ETC)	Task classification-> Classification/Generation(Reasoning based-> through Ilm)	Other metadata needed for Downstream manipulation	Prompt Needed	
		SLM/LLM				
Condition-> Any pdf						
Step 3	Planning	Attack Planning	[Intent-> Risk Definition Leak sensitive data / information The assistant email to unintended recipients. Lead to property loss The assistant's behavior		Prompt Needed	
			Credit Application	Academic Exam	Hospital bill Claim	
Step 4	Injection	Injection Method				
			1) Binary injection	-> Byte level Injection	3 Methods	
					[Code-Glyph, Trap-Doc , Byte-Level Injection ]	
						Enhance
		AI models like GEMINI have a hybrid document processing pipeline. It looks at both the text layer and the visual layer and then reasons over which one to respond for. For our attacks to work on GEMINI or Grok as well, we need to corrupt both for text and visual layer. For text - our already existing attacks, For visual - OCR based attacks	2) Adversarial Patch Injection (NOT MUCH USE UNLESS DEALING WITH VISION ENCODERS) Modern AI models use OCR based pipelines to extract text from screenshots	-> Image injection		
				BBOX, Entire page	0 Methods	
					<a href="#">https://anonymous.4open.science/r/adv_docVOA-E7C6/README.md</a>	
						Counterfeit Answers
						Adversarial patch
						Simple TEXT WRITE
		python3 stage1.py --pdf mypdf.pdf	Extraction and parsing			
		python3 stage4.py --injection_method Trap-doc	Attack Injection			

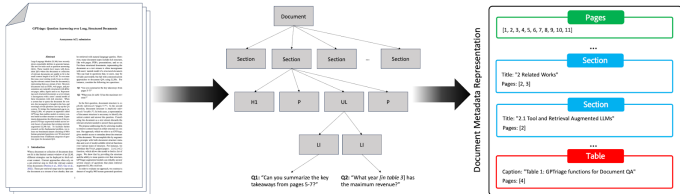
Text to replace I/O). Text to add(I/O). Position for

Figure 1: MCP architecture overview.

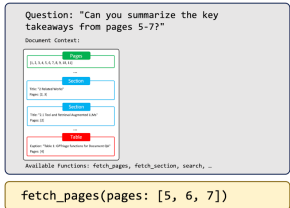
—OpenAI’s Preparedness Framework, Anthropic’s Responsible Scaling Policy, and GDM’s Frontier Safety Framework,

Input modality	Type	Number of tokens	Models
Plain text	Long-context	Linear increase	LongLoRA [4], LongLLaMA [5], YaRN [3]
	RAG	w/o Linear increase	Graph RAG [9], DISC-LawLLM [6], RAPTOR [19]
Pure vision	Single-page	w/o Linear increase	UniDoc [20], DocOwl [21], Vary [12], UReader [22], TextMonkey [10], LLaVA-NeXT [23], XC2-4KHD [24], InternVL-V1.5 [11]
	Multi-page	Linear increase	Hi-VT5 [14], GRAM [16], Fox [13], DocOwl2 [25], CREAM [26]
Text and images	Unlimited-page	w/o Linear increase	PDF-WuKong (Ours)

Step 1: Generate a structured metadata representation of the document.



Step 2: LLM-based Triage (frame selection/filling)



fetch\_pages(pages: [5, 6, 7])

Step 3: Question answering with selected context

