			Exact	Exact Contamination				
Data Contamination				Synta	Syntactic Contamination			
		Data Contamination		Exam	Examples of Each Contamination			
	Background			Signi	ficance of Conta	mination		
		Contamination from LLM Training						
		LLM Benchmarking						
		Problem Formulation						
			Math	GSM8K (Cobbe et al., 2021), MATH (Hendrycks et al., 2021), AIME 2024 (of America, 2024) and CNMO 2024 (Society, 2024).				
				ControlBench (Darioush et al., 2024), FRAMES (Krishna et al., 2024), and GPQA Diamond (Rein et al., 2023), AlpacaEval (Li et al., 2023c), ArenaHard (Li et al., 2024a).				
			Coding	(Jimene			stin et al., 2021), SWE-Bench deforces (Code-forces, 2025),	
	Static	Static Benchmarking Application	$C_{\text{Fval}}$ (Huang et al. 2024)					
	Benchmarking	1 Ippiroucion	Reasoning		(Zellers et al., 201 ARC (Clark et al.	9), WinoGrande, 2018), OpenBo	Sap et al., 2019), HellaSwag e (Sakaguchi et al., 2021), ookQA (Mihaylov et al., et al., 2018), C-SimpleQA	
			Safety			• `	al., 2020), ToxiGen	
LLM Benchmarking			Language		GLUE (Wang, 20 et al., 2019), CLU al., 2022).		E (Wang 20), Typo-fixing (Suzgun et	
			Reading Comprehens	ion	SQuAD (Rajpurk BoolQ (Clark et a		QuAC (Choi et al., 2018),	
			Canary Strin	g	BIG-Bench(Jacovi et al., 2023).			
			Encryption		Jacovi et al. (2023) TRUCE (Chandra		023),	
		Methods for Mitigation	Label Protec	tion	GLUE (Wang, 20 HumanEval (Che		E (Wang et al., 2019),	
			Direct overlap detection(Touvron et al., 2023), Radford et al., 2019; Brown et al., 2020; Chowdhery et al., 2023, Riddell et al., 2024; Lee et al., 2023; Gunasekar et al., 2023, Li et al., 2024d; Xu et al., 2024), memorization through masked inputs (Ranaldi et al., 2024; Chang et al., 2023), partial completions (Anil et al., 2023; Golchin and Surdeanu, 2024), preference for original over paraphrased test cases (Duarte et al., 2024;					
					Golchin and Surd	•		
	Problem Formulation  Correctnes  Scalability  Collision							
				Correctness				
				Complexity				
		Evaluation Criteria						
			Diversity	Diversity				
	Dynamic Benchmarking		Interpretabili	ty	Livo Donoh (W/hit	2024) A	ntil ook Donah (Wu ot ol	
			Temporal Cu	LiveBench (White et al., 2024), AntiLeak-Bench (Wu et al., 2024), AcademicEval (Zhang et al., 2024a), Live-CodeBench (Jain et al., 2024), LiveAoPSBench (Mahdavi et al., 2025), Forecastbench (Karger et al., 2024).  GSM-Symbolic (Mirzadeh et al., 2025),				
				Te	emplate-Based		(Kurtic et al., 2024), Mathador et al., 2023), MMLU-CF (Zhao	
			Rule-Based	Ta	able-Based	S3Eval (Lei et	al., 2024).	
				Gı	raph-Based	DyVal (Zhu et et al., 2024).	al., 2024a), NPHardEval (Fan	
		Existing Work			enchmark ewriting	Auto-Dataset (Cao et al., 20	(Ying et al., 2024), StructEval 24), ITD (Zhu et al., 2024c), an et al., 2024).	
			LLM-Based	HH	teractive	Interviewer (K	er (Li et al., 2023b), LM-as-an- Kim et al., 2024), TreeEval (Li KIEval (Yu et al., 2024).	
				Ev	fulti-Agent valuation	BENCHAGE	(Wang et al., 2024a), NTS (Butt et al., 2024).	
			Hybrid	l ⊢i	testEval (Li et al., 20 i et al., 2024c).	)23d), DARG (Z	Zhang et al., 2024b), C2LEVA	