

Table 1: Number of instances with compiler errors in unconstrained (Standard), idealized syntax-only (Syntax), and our proposed type-aware (Types) constraining. Type-aware constraining reduces compiler errors by 75.3% and 52.1% in the synthesis of HumanEval and MBPP problems respectively, compared to only 9.0% and 4.9% ideal improvement through syntax-only constraining on the two datasets respectively.

	Model	Synthesis			Translation		
		Standard	Syntax	Types	Standard	Syntax	Types
HumanEval	Gemma 2 2B	103	92 \downarrow 10.7%	44 \downarrow 57.3%	177	149 \downarrow 15.8%	80 \downarrow 54.8%
	Gemma 2 9B	45	41 \downarrow 8.9%	13 \downarrow 71.1%	75	63 \downarrow 16.0%	16 \downarrow 78.7%
	Gemma 2 27B	15	13 \downarrow 13.3%	2 \downarrow 86.7%	20	20 \downarrow 0.0%	3 \downarrow 85.0%
	DeepSeek Coder 33B	26	25 \downarrow 3.8%	5 \downarrow 80.8%	18	17 \downarrow 5.6%	7 \downarrow 61.1%
	CodeLlama 34B	86	71 \downarrow 17.4%	28 \downarrow 67.4%	158	124 \downarrow 21.5%	59 \downarrow 62.7%
	Qwen2.5 32B	17	17 \downarrow 0.0%	2 \downarrow 88.2%	24	21 \downarrow 12.5%	5 \downarrow 79.2%
MBPP	Gemma 2 2B	67	64 \downarrow 4.5%	27 \downarrow 59.7%	126	111 \downarrow 11.9%	79 \downarrow 37.3%
	Gemma 2 9B	30	29 \downarrow 3.3%	10 \downarrow 66.7%	67	61 \downarrow 9.0%	33 \downarrow 50.7%
	Gemma 2 27B	20	19 \downarrow 5.0%	7 \downarrow 65.0%	37	36 \downarrow 2.7%	22 \downarrow 40.5%
	DeepSeek Coder 33B	32	32 \downarrow 0.0%	19 \downarrow 40.6%	29	27 \downarrow 6.9%	13 \downarrow 55.2%
	CodeLlama 34B	80	71 \downarrow 11.2%	41 \downarrow 48.8%	126	114 \downarrow 9.5%	54 \downarrow 57.1%
	Qwen2.5 32B	19	18 \downarrow 5.3%	13 \downarrow 31.6%	22	22 \downarrow 0.0%	16 \downarrow 27.3%

Table 2: Trying to repair non-compiling instances generated with unconstrained synthesis using unconstrained decoding (Standard) and type-aware constraining (Types). The results are the number of non-repaired instances. Constrained generation boosts repair by on average 50.1%.

(a) Repair of HumanEval			(b) Repair of MBPP		
Model	Standard	Types	Model	Standard	Types
Gemma 2 2B	194 \downarrow 33.6%	103 \downarrow 64.7%	Gemma 2 2B	215 \downarrow 20.7%	125 \downarrow 53.9%
Gemma 2 9B	113 \downarrow 61.3%	52 \downarrow 82.2%	Gemma 2 9B	151 \downarrow 44.3%	82 \downarrow 69.7%
Gemma 2 27B	45 \downarrow 84.6%	22 \downarrow 92.5%	Gemma 2 27B	91 \downarrow 66.4%	50 \downarrow 81.5%
DeepSeek Coder 33B	36 \downarrow 87.7%	15 \downarrow 94.9%	DeepSeek Coder 33B	107 \downarrow 60.5%	59 \downarrow 78.2%
CodeLlama 34B	153 \downarrow 47.6%	48 \downarrow 83.6%	CodeLlama 34B	178 \downarrow 34.3%	95 \downarrow 64.9%
Qwen2.5 32B	36 \downarrow 87.7%	13 \downarrow 95.5%	Qwen2.5 32B	72 \downarrow 73.4%	45 \downarrow 83.4%

Table 3: Pass@1 (in %) of unconstrained (Standard) and type-aware constrained (Types) generated code for the tasks of Synthesis, Repair, and Translation.

	Model	Synthesis		Translation		Repair	
		Standard	Types	Standard	Types	Standard	Types
HumanEval	Gemma 2 2B	29.1	30.2	50.2	53.9	11.6	20.9
	Gemma 2 9B	56.6	58.3	73.7	78.3	24.0	34.9
	Gemma 2 27B	69.5	71.2	86.6	87.7	38.4	41.1
	DS Coder 33B	68.9	71.1	88.7	90.1	47.6	50.7
	CodeLlama 34B	41.0	43.4	58.6	63.5	17.5	27.4
	Qwen2.5 32B	79.6	81.8	92.1	93.9	65.4	71.2
MBPP	Gemma 2 2B	40.4	42.4	52.3	56.0	11.1	20.7
	Gemma 2 9B	65.4	67.4	71.4	75.8	22.1	29.2
	Gemma 2 27B	70.6	72.1	83.1	84.4	35.8	41.3
	DS Coder 33B	65.4	67.2	85.9	89.1	32.1	39.5
	CodeLlama 34B	42.2	45.6	55.7	63.3	14.8	24.7
	Qwen2.5 32B	76.3	76.6	89.6	90.4	44.6	50.2

Model	HumanEval	MBPP
Gemma 2 2B	6.7 \uparrow 38.3%	6.3 \uparrow 35.4%
Gemma 2 9B	8.3 \uparrow 29.2%	9.5 \uparrow 46.8%
Gemma 2 27B	11.7 \uparrow 19.9%	11.7 \uparrow 32.8%
DeepSeek Coder 33B	11.5 \uparrow 36.2%	9.4 \uparrow 59.5%
CodeLlama 34B	7.6 \uparrow 40.8%	7.0 \uparrow 37.6%
Qwen2.5 32B	7.3 \uparrow 39.6%	4.9 \uparrow 54.8%

Table 4: Median time taken per instance in seconds.