

APPENDIX

A. Results on Latest Close-Source Model

To further demonstrate the effectiveness and generalizability of LongCodeZip, we conduct experiments on the latest close-sourced model, Claude 3.7 Sonnet⁵. For the compression stage, we use Qwen2.5-Coder-7B as the compression model. The following tables present the results across three different tasks. We can observe that LongCodeZip still achieves significant performance improvements over other compression baselines on this powerful close-sourced model, revealing its effectiveness and generalizability.

Table VIII presents the results on the Long Code Completion task with Claude 3.7 Sonnet. Our approach demonstrates strong performance, achieving an ES score of 66.27 and an EM score of 40.20, which is on par with the no-compression baseline while operating at a 4.3x compression ratio. Notably, LongCodeZip surpasses the best-performing baseline, RAG (Function Chunking), which scored 63.55 in ES at a less efficient 3.1x ratio. This result underscores the effectiveness of our method in preserving critical information for code completion with significantly greater compression efficiency, even on state-of-the-art closed-source models.

TABLE VIII: Results on Long Code Completion with Claude 3.7 Sonnet

Model	Method	ES	EM	Ratio
CLAUDE 3.7 SONNET	<i>No Compression</i>	66.24	41.20	1.0x
	<i>No Context</i>	43.97	14.20	-
	<i>Random Token</i>	47.61	14.00	4.4x
	<i>Random Line</i>	52.61	22.20	4.5x
	RAG (Sliding Window)	61.44	34.00	2.8x
	RAG (Function Chunking)	63.55	36.80	3.1x
	LLMLingua	46.58	15.20	3.4x
	LLMLingua-2	49.02	16.20	4.4x
	LongLLMLingua	57.58	27.80	3.2x
	DietCode	54.00	19.80	3.4x
	SlimCode	53.03	20.80	4.5x
	LongCodeZip	66.27	40.20	4.3x

The results for the Long Module Summarization task are shown in Table IX. LongCodeZip remains the most competitive method, achieving a *CompScore* of 61.47, slightly outperforming the no-compression baseline (60.72) at a 1.7x compression ratio. In contrast to other baselines like RAG (58.03) and LLMLingua-2 (57.85), our approach demonstrates a clear advantage. This highlights the effectiveness of our method in preserving the most relevant semantic content for summarization, indicating that removing distracting information can even improve performance on complex summarization tasks.

⁵<https://www.anthropic.com/news/claude-3-7-sonnet>

TABLE IX: Results on Long Module Summarization with Claude 3.7 Sonnet

Model	Method	CompScore	Ratio
CLAUDE 3.7 SONNET	<i>No Compression</i>	60.72	1.0x
	<i>No Context</i>	6.58	-
	<i>Random Token</i>	37.45	1.8x
	<i>Random Line</i>	50.12	1.8x
	RAG (Sliding Window)	58.03	1.7x
	RAG (Function Chunking)	44.56	2.1x
	LLMLingua	43.21	1.7x
	LongLLMLingua	50.86	1.5x
	LLMLingua-2	57.85	2.1x
	DietCode	38.82	2.1x
	SlimCode	48.11	2.2x
	LongCodeZip	61.47	1.7x

Table X details the results on the multilingual RepoQA task. Our approach consistently achieves the best performance, with an average accuracy of 90.7% that surpasses even the no-compression baseline (89.5%) while compressing the context by 4.5x. The performance gain is particularly pronounced when compared to other compression methods; for instance, LongCodeZip outperforms the next-best baseline, LongLLMLingua (74.2%), by a margin of over 16 percentage points. This superior performance across all evaluated languages underscores our method’s effectiveness in retaining precise, structured code information necessary for retrieval-based QA.

TABLE X: Results on RepoQA with Claude 3.7 Sonnet

Method	Py	C++	Java	TS	Rust	Go	Avg. Ratio
CLAUDE 3.7 SONNET							
<i>No Compression</i>	87.4	80.1	92.6	96.7	86.3	93.6	89.5 1.0x
<i>No Context</i>	0.0	0.0	0.0	0.0	0.0	0.0	0.0 -
<i>Random Token</i>	3.6	3.1	4.2	2.1	4.2	7.3	4.1 3.6x
<i>Random Line</i>	6.2	11.4	22.9	10.4	9.4	13.5	12.3 3.5x
RAG (Sliding Window)	66.6	67.6	70.7	74.9	59.3	82.2	70.2 3.7x
RAG (Function Chunking)	56.2	48.9	61.4	40.6	60.3	71.8	56.5 4.3x
LLMLingua	5.2	7.3	9.4	11.4	4.2	16.6	9.0 4.1x
LLMLingua-2	1.0	2.1	8.3	1.0	1.0	4.2	2.9 4.6x
LongLLMLingua	72.8	65.5	73.8	70.7	81.1	81.1	74.2 4.3x
DietCode	17.7	-	36.4	-	-	-	27.1 3.7x
SlimCode	20.8	-	49.9	-	-	-	35.4 4.3x
LongCodeZip	95.7	81.1	90.5	88.4	89.4	98.8	90.7 4.5x

We also evaluated LongCodeZip on GPT-4o⁶ to further confirm its generalizability on another powerful, widely-used closed-source model. Similar to the Claude experiments, we used Qwen2.5-Coder-7B as the compressor. The results across all three tasks, presented in the tables below, affirm that our method consistently delivers top-tier performance, often matching or exceeding the no-compression baseline while operating at a high compression ratio.

Table XI shows the Long Code Completion results. LongCodeZip achieves an ES score of 64.72 and an EM score of 38.80, which is nearly identical to the no-compression baseline while compressing the context by 4.3x. This performance is notably higher than the best RAG baseline, which scored 62.01 in ES at a 3.1x ratio, highlighting our method’s superior efficiency and effectiveness on GPT-4o.

⁶<https://openai.com/index/hello-gpt-4o/>

TABLE XI: Results on Long Code Completion with GPT-4o

Model	Method	ES	EM	Ratio
GPT-4o	<i>No Compression</i>	65.13	40.80	1.0x
	<i>No Context</i>	42.92	14.00	-
	<i>Random Token</i>	46.51	13.80	4.4x
	<i>Random Line</i>	51.42	21.80	4.5x
	RAG (Sliding Window)	60.03	33.20	2.8x
	RAG (Function Chunking)	62.01	36.00	3.1x
	LLMLingua	45.53	14.80	3.4x
	LLMLingua-2	47.90	15.80	4.4x
	LongLLMLingua	56.24	27.20	3.2x
	DietCode	52.76	19.40	3.4x
	SlimCode	51.78	20.40	4.5x
	LongCodeZip	64.72	38.80	4.3x

On the Long Module Summarization task, as shown in Table [XII](#), LongCodeZip achieves a *CompScore* of 59.04, which is slightly higher than the no-compression baseline of 58.42 with a 1.7x compression. This again demonstrates that by removing noisy context, our method can help the model focus and generate better summaries.

TABLE XII: Results on Long Module Summarization with GPT-4o

Model	Method	CompScore	Ratio
GPT-4o	<i>No Compression</i>	58.42	1.0x
	<i>No Context</i>	6.41	-
	<i>Random Token</i>	35.83	1.8x
	<i>Random Line</i>	48.24	1.8x
	RAG (Sliding Window)	55.85	1.7x
	RAG (Function Chunking)	42.76	2.1x
	LLMLingua	41.57	1.7x
	LongLLMLingua	48.89	1.5x
	LLMLingua-2	55.48	2.1x
	DietCode	37.21	2.1x
	SlimCode	46.13	2.2x
	LongCodeZip	59.04	1.7x

Finally, Table [XIII](#) shows the results for the RepoQA task. LongCodeZip obtains an average accuracy of 88.9%, surpassing the no-compression baseline’s 87.8% while achieving a 4.5x compression ratio. Its performance is substantially better than the next best compression method, LongLLMLingua (72.8%). These results strongly validate the effectiveness and robustness of our approach on GPT-4o.

TABLE XIII: Results on RepoQA with GPT-4o

Method	Py	C++	Java	TS	Rust	Go	Avg. Ratio
GPT-4o							
<i>No Compression</i>	85.7	78.6	90.8	94.9	84.7	91.8	87.8 1.0x
<i>No Context</i>	0.0	0.0	0.0	0.0	0.0	0.0	0.0 -
<i>Random Token</i>	2.3	3.1	4.1	2.1	4.1	7.2	3.8 3.6x
<i>Random Line</i>	6.1	11.2	22.5	10.2	9.2	13.3	12.1 3.5x
RAG (Sliding Window)	65.3	66.3	69.4	73.5	58.2	80.6	68.9 3.7x
RAG (Function Chunking)	55.1	48.0	60.2	39.8	59.2	70.4	55.5 4.3x
LLMLingua	5.1	7.2	9.2	11.2	4.1	16.3	8.9 4.1x
LLMLingua-2	1.0	2.1	8.2	1.0	1.0	4.1	2.9 4.6x
LongLLMLingua	71.4	64.3	72.4	69.4	79.6	79.6	72.8 4.3x
DietCode	17.4	-	35.7	-	-	-	26.6 3.7x
SlimCode	20.4	-	49.0	-	-	-	34.7 4.3x
LongCodeZip	93.9	79.6	88.8	86.7	87.7	96.9	88.9 4.5x