

Table 13: Comparison among Qwen3-32B (Thinking) and other reasoning baselines. The highest and second-best scores are shown in bold and underlined, respectively.

		DeepSeek-R1 -Distill-Llama-70B	QwQ-32B	OpenAI-o3-mini (medium)	Qwen3-32B
	Architecture	Dense	Dense	-	Dense
	# Activated Params	70B	32B	-	32B
	# Total Params	70B	32B	-	32B
<i>General Tasks</i>	MMLU-Redux	89.3	<u>90.0</u>	<u>90.0</u>	<b>90.9</b>
	GPQA-Diamond	65.2	65.6	<b>76.8</b>	<u>68.4</u>
	C-Eval	71.8	<b>88.4</b>	75.1	<u>87.3</u>
	LiveBench 2024-11-25	54.5	<u>72.0</u>	70.0	<b>74.9</b>
<i>Alignment Tasks</i>	IFEval strict prompt	79.3	83.9	<b>91.5</b>	<u>85.0</u>
	Arena-Hard	60.6	<u>89.5</u>	89.0	<b>93.8</b>
	AlignBench v1.1	6.74	<u>8.70</u>	8.38	<b>8.72</b>
	Creative Writing v3	62.1	<b>82.4</b>	74.8	<u>81.0</u>
	WritingBench	6.08	<u>7.86</u>	7.52	<b>7.90</b>
<i>Math &amp; Text Reasoning</i>	MATH-500	94.5	<b>98.0</b>	<b>98.0</b>	<u>97.2</u>
	AIME'24	70.0	79.5	<u>79.6</u>	<b>81.4</b>
	AIME'25	56.3	69.5	<b>74.8</b>	<u>72.9</u>
	ZebraLogic	71.3	76.8	<b>88.9</b>	<u>88.8</u>
	AutoLogi	83.5	<b>88.1</b>	86.3	<u>87.3</u>
<i>Agent &amp; Coding</i>	BFCL v3	49.3	<u>66.4</u>	64.6	<b>70.3</b>
	LiveCodeBench v5	54.5	<u>62.7</u>	<b>66.3</b>	<u>65.7</u>
	CodeForces (Rating / Percentile)	1633 / 91.4%	<u>1982 / 97.7%</u>	<b>2036 / 98.1%</b>	1977 / 97.7%
<i>Multilingual Tasks</i>	Multi-IF	57.6	<u>68.3</u>	48.4	<b>73.0</b>
	INCLUDE	62.1	69.7	<u>73.1</u>	<b>73.7</b>
	MMMLU 14 languages	69.6	<b>80.9</b>	79.3	<u>80.6</u>
	MT-AIME2024	29.3	68.0	<u>73.9</u>	<b>75.0</b>
	PolyMath	29.4	<u>45.9</u>	38.6	<b>47.4</b>
	MLogiQA	60.3	<u>75.5</u>	71.1	<b>76.3</b>

Table 14: Comparison among Qwen3-32B (Non-thinking) and other non-reasoning baselines. The highest and second-best scores are shown in bold and underlined, respectively.

		GPT-4o-mini -2024-07-18	LLaMA-4 -Scout	Qwen2.5-72B -Instruct	Qwen3-32B
	Architecture	-	MoE	Dense	Dense
	# Activated Params	-	17B	72B	32B
	# Total Params	-	109B	72B	32B
<i>General Tasks</i>	MMLU-Redux	81.5	<u>86.3</u>	<b>86.8</b>	85.7
	GPQA-Diamond	40.2	<b>57.2</b>	49.0	<u>54.6</u>
	C-Eval	66.3	78.2	<b>84.7</b>	83.3
	LiveBench 2024-11-25	41.3	47.6	<u>51.4</u>	<b>59.8</b>
<i>Alignment Tasks</i>	IFEval strict prompt	80.4	<b>84.7</b>	<u>84.1</u>	83.2
	Arena-Hard	74.9	70.5	<u>81.2</u>	<b>92.8</b>
	AlignBench v1.1	7.81	7.49	<u>7.89</u>	<b>8.58</b>
	Creative Writing v3	<u>70.3</u>	55.0	61.8	<b>78.3</b>
	WritingBench	5.98	5.49	<u>7.06</u>	<b>7.54</b>
<i>Math &amp; Text Reasoning</i>	MATH-500	78.2	82.6	<u>83.6</u>	<b>88.6</b>
	AIME'24	8.1	<u>28.6</u>	18.9	<b>31.0</b>
	AIME'25	8.8	10.0	<u>15.0</u>	<b>20.2</b>
	ZebraLogic	20.1	24.2	<u>26.6</u>	<b>29.2</b>
	AutoLogi	52.6	56.8	<u>66.1</u>	<b>78.5</b>
<i>Agent &amp; Coding</i>	BFCL v3	<b>64.0</b>	45.4	<u>63.4</u>	63.0
	LiveCodeBench v5	27.9	29.8	<u>30.7</u>	<b>31.3</b>
	CodeForces (Rating / Percentile)	1113 / 52.6%	981 / 43.7%	859 / 35.0%	<b>1353 / 71.0%</b>
<i>Multilingual Tasks</i>	Multi-IF	62.4	64.2	<u>65.3</u>	<b>70.7</b>
	INCLUDE	66.0	<b>74.1</b>	69.6	<b>70.9</b>
	MMMLU 14 languages	72.1	<b>77.5</b>	<u>76.9</u>	76.5
	MT-AIME2024	6.0	<u>19.1</u>	12.7	<b>24.1</b>
	PolyMath	12.0	<u>20.9</u>	16.9	<b>22.5</b>
	MLogiQA	42.6	53.9	<u>59.3</u>	<b>62.9</b>