*Table 1.* **Benchmark Results of Pre-trained LLMs.** [*] indicates that LLaDA 8B Base, LLaMA2 7B Base, and LLaMA3 8B Base are evaluated under the same protocol, detailed in Appendix B.5. Results indicated by [†] and [¶] are sourced from Chu et al. (2024); Yang et al. (2024) and Bi et al. (2024) respectively. The numbers in parentheses represent the number of shots used for evaluation. "-" indicates unknown data.

| | LLaDA 8B[*] | LLaMA3 8B[*] | LLaMA2 7B[*] | Qwen2 7B[†] | Qwen2.5 7B[†] | Mistral 7B[†] | Deepseek 7B[¶] |
|---|---|---|---|---|---|---|---|
| Model | Diffusion | AR | AR | AR | AR | AR | AR |
| Training tokens | 2.3T | 15T | 2T | 7T | 18T | - | 2T |
| *General Tasks* | | | | | | | |
| MMLU | **65.9** (5) | 65.4 (5) | 45.9 (5) | 70.3 (5) | 74.2 (5) | 64.2 (5) | 48.2 (5) |
| BBH | 49.8 (3) | **57.6** (3) | 37.3 (3) | 62.3 (3) | 70.4 (3) | 56.1 (3) | 39.5 (3) |
| ARC-C | 47.9 (0) | **53.1** (0) | 46.3 (0) | 60.6 (25) | 63.7 (25) | 60.0 (25) | 48.1 (0) |
| Hellaswag | 72.5 (0) | **79.1** (0) | 76.0 (0) | 80.7 (10) | 80.2 (10) | 83.3 (10) | 75.4 (0) |
| TruthfulQA | **46.4** (0) | 44.0 (0) | 39.0 (0) | 54.2 (0) | 56.4 (0) | 42.2 (0) | - |
| WinoGrande | 74.8 (5) | **77.3** (5) | 72.5 (5) | 77.0 (5) | 75.9 (5) | 78.4 (5) | 70.5 (0) |
| PIQA | 74.4 (0) | **80.6** (0) | 79.1 (0) | - | - | - | 79.2 (0) |
| *Mathematics & Science* | | | | | | | |
| GSM8K | **70.7** (4) | 53.1 (4) | 14.3 (4) | 80.2 (4) | 85.4 (4) | 36.2 (4) | 17.4 (8) |
| Math | **27.3** (4) | 15.1 (4) | 3.2 (4) | 43.5 (4) | 49.8 (4) | 10.2 (4) | 6.0 (4) |
| GPQA | **26.1** (5) | 25.9 (5) | 25.7 (5) | 30.8 (5) | 36.4 (5) | 24.7 (5) | - |
| *Code* | | | | | | | |
| HumanEval | 33.5 (0) | **34.2** (0) | 12.8 (0) | 51.2 (0) | 57.9 (0) | 29.3 (0) | 26.2 (0) |
| HumanEval-FIM | **73.8** (2) | 73.3 (2) | 26.9 (2) | - | - | - | - |
| MBPP | 38.2 (4) | **47.4** (4) | 18.4 (4) | 64.2 (0) | 74.9 (0) | 51.1 (0) | 39.0 (3) |
| *Chinese* | | | | | | | |
| CMMLU | **69.9** (5) | 50.7 (5) | 32.5 (5) | 83.9 (5) | - | - | 47.2 (5) |
| C-Eval | **70.5** (5) | 51.7 (5) | 34.0 (5) | 83.2 (5) | - | - | 45.0 (5) |

*Table 2.* **Benchmark Results of Post-trained LLMs.** LLaDA only employs an SFT procedure while other models have extra reinforcement learning (RL) alignment. [*] indicates that LLaDA 8B Instruct, LLaMA2 7B Instruct, and LLaMA3 8B Instruct are evaluated under the same protocol, detailed in Appendix B.5. Results indicated by [†] and [¶] are sourced from Yang et al. (2024) and Bi et al. (2024) respectively. The numbers in parentheses represent the number of shots used for in-context learning. "-" indicates unknown data.

| | LLaDA 8B[*] | LLaMA3 8B[*] | LLaMA2 7B[*] | Qwen2 7B[†] | Qwen2.5 7B[†] | Gemma2 9B[†] | Deepseek 7B[¶] |
|---|---|---|---|---|---|---|---|
| Model | Diffusion | AR | AR | AR | AR | AR | AR |
| Training tokens | 2.3T | 15T | 2T | 7T | 18T | 8T | 2T |
| Post-training | SFT | SFT+RL | SFT+RL | SFT+RL | SFT+RL | SFT+RL | SFT+RL |
| Alignment pairs | 4.5M | - | - | 0.5M + - | 1M + 0.15M | - | 1.5M + - |
| *General Tasks* | | | | | | | |
| MMLU | 65.5 (5) | **68.4** (5) | 44.1 (5) | - | - | - | 49.4 (0) |
| MMLU-pro | 37.0 (0) | **41.9** (0) | 4.6 (0) | 44.1 (5) | 56.3 (5) | 52.1 (5) | - |
| Hellaswag | 74.6 (0) | **75.5** (0) | 51.5 (0) | - | - | - | 68.5 (-) |
| ARC-C | **88.5** (0) | 82.4 (0) | 57.3 (0) | - | - | - | 49.4 (-) |
| *Mathematics & Science* | | | | | | | |
| GSM8K | **78.6** (4) | 78.3 (4) | 29.0 (4) | 85.7 (0) | 91.6 (0) | 76.7 (0) | 63.0 (0) |
| Math | 26.6 (0) | **29.6** (0) | 3.8 (0) | 52.9 (0) | 75.5 (0) | 44.3 (0) | 15.8 (0) |
| GPQA | 31.8 (5) | **31.9** (5) | 28.4 (5) | 34.3 (0) | 36.4 (0) | 32.8 (0) | - |
| *Code* | | | | | | | |
| HumanEval | 47.6 (0) | **59.8** (0) | 16.5 (0) | 79.9 (0) | 84.8 (0) | 68.9 (0) | 48.2 (-) |
| MBPP | 34.2 (4) | **57.6** (4) | 20.6 (4) | 67.2 (0) | 79.2 (0) | 74.9 (0) | 35.2 (-) |