**Table 23** Checklist-Based Quality by Model, Sampling Strategy, and Task Category.

| Model | Sampling strategy | A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|---|---|
| Llama-3.1-8B-Instruct | Temperature (t=0.0) | 1.99 (0.09) | 4.54 (0.20) | 4.23 (0.27) | 3.09 (0.18) | 3.01 (0.13) | 3.57 (0.26) | 4.31 (0.13) | 4.52 (0.07) |
| Llama-3.1-8B-Instruct | Temperature (t=0.5) | 1.98 (0.09) | 4.59 (0.15) | 4.18 (0.25) | 3.02 (0.17) | 3.10 (0.13) | 3.70 (0.24) | 4.23 (0.13) | 4.53 (0.06) |
| Llama-3.1-8B-Instruct | Temperature (t=1.0) | 1.97 (0.08) | 4.48 (0.18) | 4.22 (0.23) | 2.93 (0.18) | 3.00 (0.12) | 3.61 (0.24) | 4.40 (0.09) | 4.50 (0.06) |
| Llama-3.1-8B-Instruct | In-Context Regeneration (General) | 2.21 (0.09) | 4.68 (0.13) | 4.12 (0.29) | 3.04 (0.17) | 2.85 (0.13) | 3.44 (0.25) | 4.37 (0.10) | 4.28 (0.08) |
| Llama-3.1-8B-Instruct | In-Context Regeneration (Task-Anchored) | 2.07 (0.09) | 4.59 (0.18) | 4.13 (0.23) | 2.62 (0.16) | 2.49 (0.10) | 2.93 (0.22) | 3.78 (0.14) | 3.97 (0.08) |
| Llama-3.1-8B-Instruct | System Prompt (General) | 2.22 (0.12) | 4.19 (0.17) | 4.05 (0.29) | 2.79 (0.19) | 2.38 (0.09) | 2.98 (0.19) | 3.95 (0.12) | 3.74 (0.09) |
| Llama-3.1-8B-Instruct | System Prompt (Task-Anchored) | 2.72 (0.15) | 4.42 (0.13) | 3.88 (0.27) | 2.51 (0.17) | 2.23 (0.08) | 2.48 (0.18) | 3.77 (0.15) | 3.61 (0.10) |
| Mistral-7B-Instruct-v0.3 | Temperature (t=0.0) | 2.58 (0.13) | 4.05 (0.25) | 3.37 (0.23) | 1.90 (0.15) | 3.22 (0.13) | 3.63 (0.24) | 3.75 (0.15) | 4.37 (0.07) |
| Mistral-7B-Instruct-v0.3 | Temperature (t=0.5) | 2.54 (0.12) | 3.97 (0.24) | 3.20 (0.27) | 1.83 (0.13) | 3.15 (0.13) | 3.62 (0.23) | 3.82 (0.14) | 4.39 (0.07) |
| Mistral-7B-Instruct-v0.3 | Temperature (t=1.0) | 2.41 (0.11) | 4.08 (0.21) | 3.29 (0.25) | 1.71 (0.13) | 3.12 (0.13) | 3.61 (0.22) | 3.86 (0.14) | 4.35 (0.07) |
| Mistral-7B-Instruct-v0.3 | In-Context Regeneration (General) | 2.64 (0.13) | 4.15 (0.13) | 3.44 (0.23) | 1.89 (0.14) | 2.95 (0.14) | 3.36 (0.23) | 3.75 (0.15) | 4.17 (0.07) |
| Mistral-7B-Instruct-v0.3 | In-Context Regeneration (Task-Anchored) | 2.57 (0.14) | 4.00 (0.25) | 3.60 (0.32) | 1.84 (0.15) | 2.70 (0.13) | 3.06 (0.22) | 3.45 (0.16) | 3.94 (0.08) |
| Mistral-7B-Instruct-v0.3 | System Prompt (General) | 2.48 (0.12) | 3.99 (0.21) | 3.07 (0.37) | 1.50 (0.10) | 2.24 (0.08) | 2.54 (0.18) | 3.45 (0.17) | 3.59 (0.09) |
| Mistral-7B-Instruct-v0.3 | System Prompt (Task-Anchored) | 2.55 (0.12) | 4.21 (0.17) | 3.13 (0.34) | 1.66 (0.13) | 2.08 (0.08) | 2.22 (0.17) | 3.12 (0.16) | 3.25 (0.11) |

**Table 24** Athene-RM-8B Reward by Model, Sampling Strategy, and Task Category.

| Model | Sampling Strategy | A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|---|---|
| gpt-4o | Temperature (t=0.0) | 0.37 (0.09) | 0.46 (0.13) | 0.78 (0.14) | 0.50 (0.12) | 0.28 (0.06) | 0.94 (0.21) | 0.81 (0.09) | 0.96 (0.08) |
| gpt-4o | Temperature (t=1.0) | 0.35 (0.08) | 0.43 (0.14) | 0.78 (0.12) | 0.47 (0.13) | 0.49 (0.06) | 0.97 (0.20) | 0.81 (0.10) | 0.98 (0.07) |
| gpt-4o | Temperature (t=2.0) | 0.34 (0.08) | 0.41 (0.12) | 0.80 (0.12) | 0.45 (0.14) | 0.56 (0.08) | 1.04 (0.21) | 0.85 (0.10) | 1.07 (0.07) |
| gpt-4o | In-Context Regeneration (General) | -0.38 (0.08) | -0.36 (0.11) | 0.04 (0.21) | 0.45 (0.11) | -0.40 (0.07) | 0.30 (0.19) | 0.65 (0.15) | -0.38 (0.10) |
| gpt-4o | In-Context Regeneration (Task-Anchored) | -0.08 (0.08) | -0.36 (0.10) | 0.05 (0.21) | 0.54 (0.12) | -0.24 (0.08) | 0.29 (0.24) | 0.61 (0.16) | -0.24 (0.10) |
| gpt-4o | System Prompt (General) | -0.05 (0.08) | -0.05 (0.16) | 0.15 (0.22) | 0.39 (0.13) | -1.12 (0.08) | -0.64 (0.18) | 0.22 (0.13) | -0.70 (0.10) |
| gpt-4o | System Prompt (Task-Anchored) | -0.19 (0.09) | -0.25 (0.19) | -0.13 (0.24) | -0.02 (0.12) | -1.18 (0.07) | -1.25 (0.20) | 0.16 (0.15) | -0.91 (0.11) |
| claude-4-sonnet | Temperature (t=0.0) | 0.42 (0.08) | 0.12 (0.20) | 0.69 (0.20) | 0.08 (0.15) | 1.10 (0.07) | 1.07 (0.16) | 0.96 (0.12) | 1.13 (0.07) |
| claude-4-sonnet | Temperature (t=0.5) | 0.43 (0.08) | 0.11 (0.21) | 0.68 (0.20) | 0.03 (0.15) | 1.08 (0.07) | 1.07 (0.17) | 0.97 (0.12) | 1.17 (0.07) |
| claude-4-sonnet | Temperature (t=1.0) | 0.41 (0.08) | 0.11 (0.21) | 0.67 (0.20) | 0.14 (0.14) | 1.09 (0.07) | 1.13 (0.15) | 0.96 (0.11) | 1.16 (0.07) |
| claude-4-sonnet | In-Context Regeneration (General) | 0.33 (0.08) | -0.06 (0.16) | 0.46 (0.16) | 0.25 (0.13) | 0.88 (0.07) | 0.89 (0.15) | 0.70 (0.11) | 0.78 (0.07) |
| claude-4-sonnet | In-Context Regeneration (Task-Anchored) | 0.26 (0.08) | -0.25 (0.19) | 0.42 (0.16) | -0.15 (0.11) | 0.76 (0.06) | 0.74 (0.15) | 0.48 (0.12) | 0.71 (0.07) |
| claude-4-sonnet | System Prompt (General) | 0.29 (0.08) | 0.20 (0.12) | 0.26 (0.21) | 0.39 (0.11) | -0.25 (0.08) | 0.28 (0.18) | 0.53 (0.12) | 0.29 (0.08) |
| claude-4-sonnet | System Prompt (Task-Anchored) | 0.28 (0.08) | 0.12 (0.15) | 0.29 (0.19) | 0.18 (0.11) | -0.35 (0.08) | -0.16 (0.16) | 0.50 (0.13) | -0.02 (0.10) |
| gemini-2.5-flash | Temperature (t=0.0) | 0.05 (0.10) | -0.30 (0.09) | 0.22 (0.17) | -0.20 (0.14) | 1.25 (0.15) | 0.65 (0.23) | 0.64 (0.14) | 0.56 (0.14) |
| gemini-2.5-flash | Temperature (t=1.0) | 0.12 (0.09) | -0.21 (0.10) | 0.37 (0.15) | -0.24 (0.15) | 1.27 (0.12) | 0.70 (0.26) | 0.67 (0.13) | 0.55 (0.13) |
| gemini-2.5-flash | Temperature (t=2.0) | 0.13 (0.09) | -0.31 (0.10) | 0.35 (0.15) | -0.26 (0.15) | 1.05 (0.13) | 0.69 (0.22) | 0.59 (0.12) | 0.54 (0.13) |
| gemini-2.5-flash | In-Context Regeneration (General) | -0.51 (0.09) | -0.64 (0.11) | 0.27 (0.14) | 0.04 (0.13) | 1.27 (0.09) | 0.71 (0.22) | 0.62 (0.14) | 0.45 (0.15) |
| gemini-2.5-flash | In-Context Regeneration (Task-Anchored) | -0.16 (0.09) | -0.62 (0.11) | 0.15 (0.15) | -0.28 (0.13) | 1.27 (0.09) | 0.56 (0.24) | 0.46 (0.15) | 0.44 (0.13) |
| gemini-2.5-flash | System Prompt (General) | 0.10 (0.08) | -0.28 (0.12) | -0.04 (0.23) | 0.24 (0.11) | -0.23 (0.07) | 0.16 (0.19) | 0.69 (0.11) | 0.26 (0.09) |
| gemini-2.5-flash | System Prompt (Task-Anchored) | -0.22 (0.08) | -0.56 (0.09) | 0.08 (0.20) | 0.12 (0.12) | 0.37 (0.09) | -0.62 (0.18) | 0.32 (0.16) | 0.02 (0.11) |

**Table 25** Athene-RM-8B Reward by Model, Sampling Strategy, and Task Category.

| Model | Sampling Strategy | A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|---|---|
| Llama-3.1-8B-Instruct | Temperature (t=0.0) | -0.84 (0.07) | 0.02 (0.10) | 0.03 (0.20) | -0.87 (0.18) | -0.91 (0.07) | -0.49 (0.27) | -0.22 (0.15) | 0.04 (0.08) |
| Llama-3.1-8B-Instruct | Temperature (t=0.5) | -0.81 (0.07) | 0.02 (0.12) | 0.03 (0.14) | -0.94 (0.17) | -0.90 (0.07) | -0.30 (0.18) | -0.16 (0.14) | 0.04 (0.07) |
| Llama-3.1-8B-Instruct | Temperature (t=1.0) | -0.95 (0.06) | -0.17 (0.16) | 0.14 (0.13) | -0.86 (0.15) | -0.76 (0.07) | -0.29 (0.17) | -0.06 (0.12) | 0.06 (0.07) |
| Llama-3.1-8B-Instruct | In-Context Regeneration (General) | -0.62 (0.06) | -0.34 (0.15) | -0.21 (0.23) | -0.92 (0.17) | -1.34 (0.08) | -0.87 (0.14) | -0.29 (0.14) | -0.51 (0.08) |
| Llama-3.1-8B-Instruct | In-Context Regeneration (Task-Anchored) | -0.74 (0.06) | -0.57 (0.12) | -0.09 (0.16) | -1.47 (0.16) | -1.59 (0.08) | -1.24 (0.13) | -0.75 (0.15) | -0.92 (0.09) |
| Llama-3.1-8B-Instruct | System Prompt (General) | -0.24 (0.08) | -0.36 (0.22) | -0.02 (0.19) | -0.83 (0.16) | -2.08 (0.09) | -1.53 (0.17) | -0.80 (0.17) | -1.35 (0.11) |
| Llama-3.1-8B-Instruct | System Prompt (Task-Anchored) | -0.14 (0.09) | -0.19 (0.19) | 0.00 (0.22) | -1.05 (0.15) | -2.24 (0.09) | -2.04 (0.19) | -0.86 (0.18) | -1.48 (0.12) |
| Mistral-7B-Instruct-v0.3 | Temperature (t=0.0) | -0.08 (0.07) | -0.08 (0.20) | -0.78 (0.15) | -1.30 (0.15) | -1.12 (0.08) | -0.45 (0.13) | -0.35 (0.13) | -0.39 (0.07) |
| Mistral-7B-Instruct-v0.3 | Temperature (t=0.5) | -0.06 (0.06) | -0.07 (0.16) | -0.84 (0.15) | -1.37 (0.13) | -1.06 (0.08) | -0.39 (0.11) | -0.29 (0.11) | -0.32 (0.06) |
| Mistral-7B-Instruct-v0.3 | Temperature (t=1.0) | -0.09 (0.06) | -0.17 (0.16) | -0.74 (0.14) | -1.40 (0.12) | -0.98 (0.07) | -0.36 (0.12) | -0.18 (0.10) | -0.28 (0.06) |
| Mistral-7B-Instruct-v0.3 | In-Context Regeneration (General) | -0.18 (0.07) | -0.33 (0.13) | -0.74 (0.20) | -1.32 (0.13) | -1.22 (0.09) | -0.76 (0.13) | -0.46 (0.12) | -0.73 (0.08) |
| Mistral-7B-Instruct-v0.3 | In-Context Regeneration (Task-Anchored) | -0.10 (0.07) | -0.48 (0.15) | -0.59 (0.21) | -1.38 (0.13) | -1.54 (0.11) | -1.01 (0.17) | -1.00 (0.12) | -1.07 (0.09) |
| Mistral-7B-Instruct-v0.3 | System Prompt (General) | -0.41 (0.08) | -0.49 (0.25) | -1.17 (0.22) | -1.81 (0.13) | -2.19 (0.08) | -2.06 (0.21) | -1.04 (0.17) | -1.81 (0.13) |
| Mistral-7B-Instruct-v0.3 | System Prompt (Task-Anchored) | -0.41 (0.08) | -0.33 (0.18) | -1.20 (0.25) | -1.86 (0.15) | -2.30 (0.09) | -2.46 (0.19) | -1.29 (0.16) | -2.15 (0.15) |