

## FOOTNOTE URLs

- <sup>2</sup><https://aws.amazon.com/blogs/machine-learning/fmops-llmops-operation-alize-generative-ai-and-differences-with-mlops>
- <sup>3</sup><https://huggingface.co/spaces>
- <sup>4</sup><https://pages.github.com>
- <sup>5</sup><https://paperswithcode.com>
- <sup>6</sup>[https://huggingface.co/spaces/open-llm-leaderboard/open\\_llm\\_leaderboard/discussions/801](https://huggingface.co/spaces/open-llm-leaderboard/open_llm_leaderboard/discussions/801)
- <sup>7</sup><https://huggingface.co/blog/open-llm-leaderboard-drop>
- <sup>8</sup><https://huggingface.co/spaces/open-llm-leaderboard/blog>
- <sup>9</sup><https://github.com/topics/awesome-list>
- <sup>10</sup><https://scholar.google.com>
- <sup>11</sup><https://github.com/SAILResearch/awesome-foundation-model-leaderboards>
- <sup>12</sup><https://huggingface.co/spaces/zhiminy/awesome-foundation-model-leaderboard-search>
- <sup>13</sup><https://chat.lmsys.org/?leaderboard>
- <sup>14</sup><https://www.datacomp.ai/dclm/leaderboard.html>
- <sup>15</sup><https://www.datacomp.ai/dclclip/leaderboard.html>
- <sup>16</sup><https://github.com/yuh-zha/AlignScore>
- <sup>17</sup>[https://huggingface.co/spaces/locuslab/tofu\\_leaderboard](https://huggingface.co/spaces/locuslab/tofu_leaderboard)
- <sup>18</sup><https://huggingface.co/spaces/ameerazam08/Paper-LeaderBoard>
- <sup>19</sup><https://leaderboard.withmartian.com>
- <sup>20</sup><https://www.kaggle.com/competitions>
- <sup>21</sup><https://eval.ai>
- <sup>22</sup><https://taostats.io>
- <sup>23</sup><https://huggingface.co/spaces/mrfakename/open-leaderboards-leaderboard>
- <sup>24</sup><https://super.gluebenchmark.com/leaderboard>
- <sup>25</sup><https://neptune.ai/blog/llmops>
- <sup>26</sup><https://github.com/paperswithcode/paperswithcode-data>
- <sup>27</sup><https://huggingface.co/spaces?sort=trending&search=leaderboard>
- <sup>28</sup>[https://huggingface.co/spaces/Vikhrmodels/Russian\\_Arena\\_Hard/discussions/3](https://huggingface.co/spaces/Vikhrmodels/Russian_Arena_Hard/discussions/3)
- <sup>29</sup><https://huggingface.co/spaces?search=chatbot+arena+leaderboard>
- <sup>30</sup><https://github.com/Helsinki-NLP/OPUS-MT-dashboard/issues/2>
- <sup>31</sup>[https://huggingface.co/spaces/HuggingFaceH4/human\\_eval\\_llm\\_leaderboard/discussions/4](https://huggingface.co/spaces/HuggingFaceH4/human_eval_llm_leaderboard/discussions/4)
- <sup>32</sup><https://www.superclueai.com>
- <sup>33</sup><https://github.com/SAILResearch/awesome-foundation-model-leaderboards>
- <sup>34</sup><https://huggingface.co/spaces/zhiminy/awesome-foundation-model-leaderboard-search>
- <sup>35</sup><https://slack.com>
- <sup>36</sup><https://discord.com>
- <sup>37</sup><https://weixin.qq.com>
- <sup>38</sup><https://github.com/THU-KEG/KoLA/issues/17>
- <sup>39</sup><https://github.com/embeddings-benchmark/mteb/pull/187>
- <sup>40</sup><https://github.com/LudwigStumpp/llm-leaderboard>
- <sup>41</sup>[https://tatsu-lab.github.io/alpaca\\_eval](https://tatsu-lab.github.io/alpaca_eval)
- <sup>42</sup><https://huggingface.co/spaces/bigcode/bigcode-models-leaderboard>
- <sup>43</sup><https://github.com/paperswithcode/axcell>
- <sup>44</sup><https://huggingface.co/spaces/gaia-benchmark/leaderboard>
- <sup>45</sup><https://leaderboard.allenai.org/a-okvqa/submissions/public>
- <sup>46</sup><https://cmedbenchmark.llmzoo.com/static/leaderboard.html>
- <sup>47</sup><https://opencompass.org.cn/doc>
- <sup>48</sup><https://jwolpxeehx.feishu.cn/wiki/C6VfwvbmOiuVrokpJAgcJXUcnLh>
- <sup>49</sup>[https://docs.google.com/document/d/1Rs6d\\_pcs2vfqW4Ymub7Wb1XtBNlrc-WfH3T7U1ktuBo](https://docs.google.com/document/d/1Rs6d_pcs2vfqW4Ymub7Wb1XtBNlrc-WfH3T7U1ktuBo)
- <sup>50</sup><https://github.com/stars>
- <sup>51</sup>[https://huggingface.co/spaces/NexaAIDev/domain\\_llm\\_leaderboard](https://huggingface.co/spaces/NexaAIDev/domain_llm_leaderboard)
- <sup>52</sup><https://flageval.baai.ac.cn/#/leaderboard>
- <sup>53</sup><https://cmmmu-benchmark.github.io/#leaderboard>
- <sup>54</sup><https://research.ibm.com/blog/efficient-llm-benchmarking>
- <sup>55</sup><https://chat.lmsys.org>
- <sup>56</sup><https://huggingface.co/spaces/Auto-Arena/Leaderboard>
- <sup>57</sup><https://rank.opencompass.org.cn/leaderboard-arena>
- <sup>58</sup><https://huggingface.co/spaces/llm-council/emotional-intelligence-arena>
- <sup>59</sup><https://bryanyzhu.github.io/posts/2024-06-20-elo-part1>
- <sup>60</sup><https://arena.lmsys.org>
- <sup>61</sup>[https://twitter.com/eval\\_ai/status/1341091250544521231](https://twitter.com/eval_ai/status/1341091250544521231)
- <sup>62</sup><https://github.com/paperswithcode/model-index>
- <sup>63</sup>[https://github.com/rowanzh/hellaswag/tree/master/hellaswag\\_models#submitting-to-the-leaderboard](https://github.com/rowanzh/hellaswag/tree/master/hellaswag_models#submitting-to-the-leaderboard)
- <sup>64</sup>[https://github.com/tatsu-lab/alpaca\\_eval?tab=readme-ov-file#contributing-a-model](https://github.com/tatsu-lab/alpaca_eval?tab=readme-ov-file#contributing-a-model)
- <sup>65</sup><https://github.com/ray-project/llmperf-leaderboard?tab=readme-ov-file#feedback>
- <sup>66</sup>[https://huggingface.co/spaces/HuggingFaceH4/open\\_llm\\_leaderboard](https://huggingface.co/spaces/HuggingFaceH4/open_llm_leaderboard)
- <sup>67</sup>[https://huggingface.co/spaces/Cognitive-Lab/indic\\_llm\\_leaderboard](https://huggingface.co/spaces/Cognitive-Lab/indic_llm_leaderboard)
- <sup>68</sup><https://crfm.stanford.edu/helm/lite/latest/#/leaderboard>
- <sup>69</sup><https://stats.stackexchange.com/questions/333446/what-does-the-term-gold-label-refer-to-in-the-context-of-semi-supervised-class>
- <sup>70</sup><https://github.com/EleutherAI/lm-evaluation-harness>
- <sup>71</sup><https://github.com/lm-sys/FastChat>
- <sup>72</sup>[https://tatsu-lab.github.io/alpaca\\_eval](https://tatsu-lab.github.io/alpaca_eval)
- <sup>73</sup>[https://huggingface.co/spaces/gsaivinay/open\\_llm\\_leaderboard](https://huggingface.co/spaces/gsaivinay/open_llm_leaderboard)
- <sup>74</sup><https://github.com/MikeGu721/XiezhiBenchmark>
- <sup>75</sup><https://leaderboard.tabbyml.com>
- <sup>76</sup><https://artificialanalysis.ai/leaderboards/models>
- <sup>77</sup><https://videoniah.github.io>
- <sup>78</sup><https://xwang.dev/mint-bench>
- <sup>79</sup><https://osu-nlp-group.github.io/TravelPlanner>
- <sup>80</sup>[https://huggingface.co/spaces/BramVanroy/open\\_dutch\\_llm\\_leaderboard](https://huggingface.co/spaces/BramVanroy/open_dutch_llm_leaderboard)
- <sup>81</sup><https://huggingface.co/spaces/ml-energy/leaderboard>
- <sup>82</sup>[https://huggingface.co/spaces/ramirolo/LLMHallucination\\_Leaderboard](https://huggingface.co/spaces/ramirolo/LLMHallucination_Leaderboard)
- <sup>83</sup><https://huggingface.co/spaces/dimbyTa/open-llm-leaderboard-viz>
- <sup>84</sup><https://huggingface.co/spaces/upstage/open-ko-llm-leaderboard/discussions/28>
- <sup>85</sup><https://github.com/GAIR-NLP/factool/issues/39>
- <sup>86</sup><https://github.com/microsoft/promptbench/issues/26>
- <sup>87</sup><https://github.com/nyu-ml/jiant/issues/1366>
- <sup>88</sup><https://github.com/TheoremOne/llm-benchmark-suite/issues/10>
- <sup>89</sup><https://huggingface.co/spaces/lovodkin93/FuseReviews-Leaderboard/discussions/1>
- <sup>90</sup><https://huggingface.co/spaces/nlphuji/WHOOOPS-Leaderboard-Full/discussions/2>
- <sup>91</sup>[https://huggingface.co/spaces/Nexusflow/Nexus\\_Function\\_Calling\\_Leaderboard/discussions/3](https://huggingface.co/spaces/Nexusflow/Nexus_Function_Calling_Leaderboard/discussions/3)
- <sup>92</sup>[https://huggingface.co/spaces/SeaEval/SeaEval\\_Leaderboard/discussions/2](https://huggingface.co/spaces/SeaEval/SeaEval_Leaderboard/discussions/2)
- <sup>93</sup><https://github.com/open-compass/LawBench/issues/7>
- <sup>94</sup><https://huggingface.co/spaces/mike-ravkine/can-ai-code-results/discussions/4>
- <sup>95</sup><https://huggingface.co/spaces/mike-ravkine/can-ai-code-results/discussions/6>
- <sup>96</sup><https://huggingface.co/spaces/mike-ravkine/can-ai-code-results/discussions/9>
- <sup>97</sup><https://github.com/tjunlp-lab/M3KE/issues/8>
- <sup>98</sup><https://github.com/open-compass/MathBench/issues/7>
- <sup>99</sup><https://github.com/BradyFU/Awesome-Multimodal-Large-Language-Models/issues/124>
- <sup>100</sup><https://github.com/hendrycks/test/issues/25>
- <sup>101</sup><https://github.com/princeton-nlp/SWE-bench/issues/46>
- <sup>102</sup><https://github.com/CLUEbenchmark/SuperCLUE/issues/43>
- <sup>103</sup><https://github.com/CLUEbenchmark/SuperCLUE-Agent/issues/8>
- <sup>104</sup><https://github.com/CLUEbenchmark/SuperCLUE-Auto/issues/3>
- <sup>105</sup><https://github.com/CLUEbenchmark/SuperCLUE-Math6/issues/1>
- <sup>106</sup><https://github.com/CLUEbenchmark/SuperCLUE-Safety/issues/8>
- <sup>107</sup><https://github.com/CLUEbenchmark/SuperCLUE-Egkzw/issues/2>
- <sup>108</sup><https://github.com/CLUEbenchmark/SuperCLUE-RAG/issues/2>
- <sup>109</sup><https://github.com/CLUEbenchmark/SuperCLUE-Code3/issues/3>
- <sup>110</sup><https://github.com/CLUEbenchmark/SuperCLUE-Role/issues/2>
- <sup>111</sup><https://github.com/CLUEbenchmark/SuperCLUE-Industry/issues/1>
- <sup>112</sup>[https://github.com/EleutherAI/lm\\_perplexity/issues/2](https://github.com/EleutherAI/lm_perplexity/issues/2)
- <sup>113</sup><https://huggingface.co/spaces/mike-ravkine/can-ai-code-results/discussions/7>
- <sup>114</sup><https://huggingface.co/spaces/mike-ravkine/can-ai-code-results/discussions/3>
- <sup>115</sup><https://github.com/stanford-crfm/helm/issues/2386>
- <sup>116</sup><https://github.com/hkust-nlp/ceval/issues/76>
- <sup>117</sup><https://huggingface.co/spaces/daisheh/SCULAiW/discussions/1>
- <sup>118</sup><https://github.com/MikeGu721/XiezhiBenchmark/issues/6>
- <sup>119</sup><https://github.com/OpenLM/Lab/LEval/issues/10>
- <sup>120</sup><https://huggingface.co/spaces/ought/raft-leaderboard/discussions/20>
- <sup>121</sup><https://github.com/taoyds/spider/issues/101>
- <sup>122</sup><https://github.com/embeddings-benchmark/mteb/issues/186>
- <sup>123</sup><https://github.com/Helsinki-NLP/OPUS-MT-dashboard/issues/2>
- <sup>124</sup>[https://huggingface.co/spaces/opencompass/open\\_vlm\\_leaderboard/discussions/6](https://huggingface.co/spaces/opencompass/open_vlm_leaderboard/discussions/6)
- <sup>125</sup><https://github.com/LaVi-Lab/CLEVA/issues/9>
- <sup>126</sup><https://github.com/FlagOpen/FlagEval/issues/36>
- <sup>127</sup><https://github.com/palash1992/GEM-Benchmark/issues/4>

<sup>128</sup><https://github.com/YangLinyi/GLUE-X/issues/4>  
<sup>129</sup><https://github.com/OpenGVLab/Multi-Modality-Arena/issues/19>  
<sup>130</sup><https://github.com/spraakbanken/SuperLim-2/issues/1>  
<sup>131</sup><https://huggingface.co/spaces/nlphuji/WHOOPS-Leaderboard-Full/discussions/1>  
<sup>132</sup><https://huggingface.co/spaces/mlfoundations/VisIT-Bench-Leaderboard/discussions/1>  
<sup>133</sup><https://github.com/paperswithcode/sota-extractor/issues/25>  
<sup>134</sup>[https://huggingface.co/spaces/mlabonne/Yet\\_Another\\_LLM\\_Leaderboard/discussions/8](https://huggingface.co/spaces/mlabonne/Yet_Another_LLM_Leaderboard/discussions/8)  
<sup>135</sup><https://github.com/Q-Future/Q-Bench/issues/11>  
<sup>136</sup><https://github.com/stanford-crfm/helm/issues/2034>  
<sup>137</sup><https://github.com/stanford-crfm/helm/issues/2060>  
<sup>138</sup><https://github.com/LudwigStumpp/llm-leaderboard/issues/9>  
<sup>139</sup><https://github.com/google/BIG-bench/issues/983>  
<sup>140</sup><https://web.archive.org/web/20221205184917/https://paperswithcode.com/sota/image-classification-on-humaneval>  
<sup>141</sup><https://github.com/LudwigStumpp/llm-leaderboard/issues/10>  
<sup>142</sup><https://github.com/FlagOpen/FlagEval/issues/34>  
<sup>143</sup><https://github.com/THU-KEG/KoLA/issues/5>  
<sup>144</sup><https://github.com/OpenM3D/M3DBench/issues/2>  
<sup>145</sup><https://github.com/OpenLMMLab/GAOKAO-Bench/issues/25>  
<sup>146</sup>[https://github.com/thodan/bop\\_toolkit/issues/119](https://github.com/thodan/bop_toolkit/issues/119)  
<sup>147</sup><https://github.com/princeton-nlp/intercode/issues/23>  
<sup>148</sup><https://github.com/stanford-crfm/helm/issues/2238>  
<sup>149</sup><https://github.com/stanford-crfm/helm/issues/2008>  
<sup>150</sup><https://github.com/stanford-crfm/helm/issues/2238>  
<sup>151</sup>[https://huggingface.co/spaces/Wwwduojin/MLLM\\_leaderboard/discussions/1](https://huggingface.co/spaces/Wwwduojin/MLLM_leaderboard/discussions/1)  
<sup>152</sup><https://paperswithcode.com/sota/zero-shot-video-question-answer-on-starr>  
<sup>153</sup><https://paperswithcode.com/sota/zero-shot-video-question-answer-on-starr-1>  
<sup>154</sup><https://github.com/THU-KEG/KoLA/issues/6>  
<sup>155</sup><https://github.com/THU-KEG/KoLA/issues/7>  
<sup>156</sup><https://github.com/THU-KEG/KoLA/issues/2>  
<sup>157</sup>[https://huggingface.co/spaces/OpenGVLab/MVBench\\_Leaderboard/discussions/4](https://huggingface.co/spaces/OpenGVLab/MVBench_Leaderboard/discussions/4)  
<sup>158</sup>[https://huggingface.co/spaces/OpenGVLab/MVBench\\_Leaderboard/discussions/3](https://huggingface.co/spaces/OpenGVLab/MVBench_Leaderboard/discussions/3)  
<sup>159</sup><https://github.com/THU-KEG/KoLA/issues/17>  
<sup>160</sup><https://github.com/stanford-crfm/helm/issues/2351>  
<sup>161</sup>[https://huggingface.co/docs/huggingface\\_hub/package\\_reference/cards](https://huggingface.co/docs/huggingface_hub/package_reference/cards)  
<sup>162</sup><https://huggingface.co/spaces/demo-leaderboard-backend/leaderboard>  
<sup>163</sup>[https://huggingface.co/spaces/freddyaboulton/gradio\\_leaderboard](https://huggingface.co/spaces/freddyaboulton/gradio_leaderboard)  
<sup>164</sup><https://hub.opencompass.org.cn/dataset-submit>  
<sup>165</sup><https://github.com/xingyaoww/mint-bench/blob/main/docs/CONTRIBUTING.md#contribute-dataset>  
<sup>166</sup>[https://github.com/tatsu-lab/alpaca\\_eval?tab=readme-ov-file#contributing-an-eval-set](https://github.com/tatsu-lab/alpaca_eval?tab=readme-ov-file#contributing-an-eval-set)  
<sup>167</sup><https://github.com/llm-council/llm-council?tab=readme-ov-file#adding-a-prompt>  
<sup>168</sup><https://github.com/svilupp/Julia-LLM-Leaderboard#contributing-your-test-case>  
<sup>169</sup>[https://github.com/tatsu-lab/alpaca\\_eval?tab=readme-ov-file#contributing-an-evaluator](https://github.com/tatsu-lab/alpaca_eval?tab=readme-ov-file#contributing-an-evaluator)  
<sup>170</sup>[https://github.com/tatsu-lab/alpaca\\_eval?tab=readme-ov-file#contributing-a-completion-function](https://github.com/tatsu-lab/alpaca_eval?tab=readme-ov-file#contributing-a-completion-function)  
<sup>171</sup><https://twitter.com/karpathy/status/1737544497016578453>  
<sup>172</sup>[https://twitter.com/Nils\\_Reimers/status/1777653735192556621](https://twitter.com/Nils_Reimers/status/1777653735192556621)  
<sup>173</sup><https://huggingface.co/spaces/andrewreed/closed-vs-open-arena-elo>