Table 3 | We use OmniDocBench [27] to test the performance of DeepSeek-OCR on real document parsing tasks. All metrics in the table are edit distances, where smaller values indicate better performance. "Tokens" represents the average number of vision tokens used per page, and "†200dpi" means using *fitz* to interpolate the original image to 200dpi. For the DeepSeek-OCR model, the values in parentheses in the "Tokens" column represent valid vision tokens, calculated according to Equation 1.

| Model | Tokens | English | | | | | Chinese | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | overall | text | formula | table | order | overall | text | formula | table | order |
| **Pipline Models** | | | | | | | | | | | |
| Dolphin [11] | - | 0.356 | 0.352 | 0.465 | 0.258 | 0.35 | 0.44 | 0.44 | 0.604 | 0.367 | 0.351 |
| Marker [1] | - | 0.296 | 0.085 | 0.374 | 0.609 | 0.116 | 0.497 | 0.293 | 0.688 | 0.678 | 0.329 |
| Mathpix [2] | - | 0.191 | 0.105 | 0.306 | 0.243 | 0.108 | 0.364 | 0.381 | 0.454 | 0.32 | 0.30 |
| MinerU-2.1.1 [34] | - | 0.162 | 0.072 | 0.313 | 0.166 | 0.097 | 0.244 | 0.111 | 0.581 | 0.15 | 0.136 |
| MonkeyOCR-1.2B [18] | - | 0.154 | 0.062 | 0.295 | 0.164 | 0.094 | 0.263 | 0.179 | 0.464 | 0.168 | 0.243 |
| PPstructure-v3 [9] | - | 0.152 | 0.073 | 0.295 | 0.162 | 0.077 | 0.223 | 0.136 | 0.535 | 0.111 | 0.11 |
| **End-to-end Models** | | | | | | | | | | | |
| Nougat [6] | 2352 | 0.452 | 0.365 | 0.488 | 0.572 | 0.382 | 0.973 | 0.998 | 0.941 | 1.00 | 0.954 |
| SmolDocling [25] | 392 | 0.493 | 0.262 | 0.753 | 0.729 | 0.227 | 0.816 | 0.838 | 0.997 | 0.907 | 0.522 |
| InternVL2-76B [8] | 6790 | 0.44 | 0.353 | 0.543 | 0.547 | 0.317 | 0.443 | 0.29 | 0.701 | 0.555 | 0.228 |
| Qwen2.5-VL-7B [5] | 3949 | 0.316 | 0.151 | 0.376 | 0.598 | 0.138 | 0.399 | 0.243 | 0.5 | 0.627 | 0.226 |
| OLMOCR [28] | 3949 | 0.326 | 0.097 | 0.455 | 0.608 | 0.145 | 0.469 | 0.293 | 0.655 | 0.652 | 0.277 |
| GOT-OCR2.0 [38] | 256 | 0.287 | 0.189 | 0.360 | 0.459 | 0.141 | 0.411 | 0.315 | 0.528 | 0.52 | 0.28 |
| OCRFlux-3B [3] | 3949 | 0.238 | 0.112 | 0.447 | 0.269 | 0.126 | 0.349 | 0.256 | 0.716 | 0.162 | 0.263 |
| GPT4o [26] | - | 0.233 | 0.144 | 0.425 | 0.234 | 0.128 | 0.399 | 0.409 | 0.606 | 0.329 | 0.251 |
| InternVL3-78B [42] | 6790 | 0.218 | 0.117 | 0.38 | 0.279 | 0.095 | 0.296 | 0.21 | 0.533 | 0.282 | 0.161 |
| Qwen2.5-VL-72B [5] | 3949 | 0.214 | 0.092 | 0.315 | 0.341 | 0.106 | 0.261 | 0.18 | 0.434 | 0.262 | 0.168 |
| dots.ocr [30] | 3949 | 0.182 | 0.137 | 0.320 | 0.166 | 0.182 | 0.261 | 0.229 | 0.468 | 0.160 | 0.261 |
| Gemini2.5-Pro [4] | - | 0.148 | 0.055 | 0.356 | 0.13 | 0.049 | 0.212 | 0.168 | 0.439 | 0.119 | 0.121 |
| MinerU2.0 [34] | 6790 | 0.133 | 0.045 | 0.273 | 0.15 | 0.066 | 0.238 | 0.115 | 0.506 | 0.209 | 0.122 |
| dots.ocr†200dpi [30] | 5545 | 0.125 | **0.032** | 0.329 | **0.099** | **0.04** | 0.16 | **0.066** | 0.416 | 0.092 | **0.067** |
| **DeepSeek-OCR (end2end)** | | | | | | | | | | | |
| Tiny | **64** | 0.386 | 0.373 | 0.469 | 0.422 | 0.283 | 0.361 | 0.307 | 0.635 | 0.266 | 0.236 |
| Small | 100 | 0.221 | 0.142 | 0.373 | 0.242 | 0.125 | 0.284 | 0.24 | 0.53 | 0.159 | 0.205 |
| Base | 256(182) | 0.137 | 0.054 | 0.267 | 0.163 | 0.064 | 0.24 | 0.205 | 0.474 | 0.1 | 0.181 |
| Large | 400(285) | 0.138 | 0.054 | 0.277 | 0.152 | 0.067 | 0.208 | 0.143 | 0.461 | 0.104 | 0.123 |
| Gundam | 795 | 0.127 | 0.043 | 0.269 | 0.134 | 0.062 | 0.181 | 0.097 | 0.432 | 0.089 | 0.103 |
| Gundam-M†200dpi | 1853 | **0.123** | 0.049 | **0.242** | 0.147 | 0.056 | **0.157** | 0.087 | **0.377** | **0.08** | 0.085 |

without layout: "<image>\nFree OCR." to control the model's output format. Nevertheless, the output format still cannot completely match Fox benchmarks, so the actual performance would be somewhat higher than the test results.

As shown in Table 2, within a 10× compression ratio, the model's decoding precision can reach approximately 97%, which is a very promising result. In the future, it may be possible to achieve nearly 10× lossless contexts compression through text-to-image approaches. When the compression ratio exceeds 10×, performance begins to decline, which may have two reasons: one is that the layout of long documents becomes more complex, and another reason may be that long texts become blurred at 512×512 or 640×640 resolution. The first issue can be solved by rendering texts onto a single layout page, while we believe the second issue will become