

# Fuentes Académicas sobre el Mercado de Conversiones de Archivos y Nuevos Formatos de Documentos

## Introducción

El campo de la conversión de documentos y formatos de archivos está experimentando una transformación significativa impulsada por la inteligencia artificial, el procesamiento de documentos inteligente y los avances en tecnologías de código abierto. Este informe presenta las fuentes académicas más valoradas, innovadoras y recientes sobre el mercado de conversiones de archivos, formatos emergentes y procesos de conversión open source.

## Fuentes Académicas Principales sobre Conversión de Documentos

### Docling: Innovación en Conversión de Documentos con IA

**Docling: An Efficient Open-Source Toolkit for AI-driven Document Conversion** <sup>[1]</sup> representa una de las contribuciones más significativas y recientes al campo. Este trabajo, publicado en enero de 2025, presenta un toolkit de código abierto bajo licencia MIT que utiliza modelos de IA especializados para análisis de layout (DocLayNet) y reconocimiento de estructura de tablas (TableFormer). La herramienta procesa formatos populares como PDF, imágenes, MS Office y HTML, convirtiéndolos a representaciones estructuradas unificadas.

#### Características técnicas destacadas:

- Ejecuta completamente en hardware local
- Arquitectura modular y extensible
- Integración nativa con LangChain, LlamaIndex y spaCy
- Procesamiento eficiente con aceleración GPU
- Más de 10,000 estrellas en GitHub en menos de un mes

### Benchmarking y Evaluación Comparativa

**A Comprehensive and Challenging OCR Benchmark for Evaluating LMMs** <sup>[2]</sup> introduce CC-OCR, un benchmark integral diseñado específicamente para evaluar modelos de lenguaje multimodal en diferentes dimensiones de capacidades de alfabetización. Este trabajo es crucial para establecer métricas de evaluación estándar en el campo.

**OmniDocBench: Benchmarking Diverse PDF Document Parsing** <sup>[3]</sup> presenta un benchmark meticulosamente curado que incluye nueve tipos diversos de documentos con anotaciones ricas, incluyendo anotaciones de layout y reconocimiento. Los resultados muestran que

herramientas pipeline especializadas como MinerU y Mathpix superan a los modelos de lenguaje general.

## Procesamiento Inteligente de Documentos (IDP)

**Document Parsing Unveiled: Techniques, Challenges, and Prospects** <sup>[4]</sup> ofrece una revisión integral del estado actual del análisis de documentos, cubriendo metodologías clave desde sistemas de pipeline modulares hasta modelos end-to-end impulsados por grandes modelos de visión-lenguaje.

**ERPA: Efficient RPA Model Integrating OCR and LLMs** <sup>[5]</sup> propone un marco RPA mejorado que integra tecnología OCR avanzada con modelos de lenguaje grandes, logrando una reducción del 93% en tiempo de procesamiento en comparación con plataformas RPA existentes.

## Nuevos Formatos de Documentos Emergentes

### Formatos de Nueva Generación

**OME-NGFF: scalable format strategies for interoperable bioimaging** <sup>[6]</sup> introduce formatos de archivo de próxima generación (NGFF) como solución a los desafíos de compatibilidad. Los formatos incluyen TIFF, HDF5 y Zarr, diseñados para proporcionar el equilibrio necesario para que la comunidad converja y reduzca el desarrollo de más formatos.

**WebAssembly Text Format** <sup>[7]</sup> representa un avance en formatos de representación textual legible por humanos. WebAssembly tiene dos tipos de archivos: `.wasm` (bytecode) y `.wat` (representación textual legible), siendo este último crucial para depuración e inspección.

### Estándares de Documentos Modernos

**Standards for language resources in ISO** <sup>[8]</sup> documenta el portafolio de estándares que cubre principios básicos para identificar recursos de lenguaje, documentación (metadatos) y gestión de ciclo de vida, fundamental para el desarrollo de nuevos formatos.

**PDF Standards: Specifications of PDF/A, PDF/UA and More** <sup>[9]</sup> detalla la evolución de estándares PDF, incluyendo PDF/A para archivado, PDF/X para intercambio gráfico, PDF/E para ingeniería, y PDF/UA para accesibilidad universal.

## Procesos de Conversión Open Source

### Herramientas de Conversión de Código Abierto

**EffOCR: An Extensible, Open-Source Package for Efficiently Digitizing World Knowledge** <sup>[10]</sup> presenta un paquete OCR de código abierto diseñado para investigadores que requiere ser preciso, extremadamente barato de implementar y eficiente en muestras para personalizar a nuevas colecciones.

**appjsonify: An Academic Paper PDF-to-JSON Conversion Toolkit** <sup>[11]</sup> presenta un toolkit basado en Python para conversión PDF-a-JSON específicamente para papers académicos,

utilizando modelos de análisis de layout basados en visión y enfoques de procesamiento de texto basados en reglas.

**OnPrem.LLM: A Privacy-Conscious Document Intelligence Toolkit** <sup>[12]</sup> ofrece un toolkit basado en Python para aplicar modelos de lenguaje grandes a datos sensibles en entornos offline o restringidos, proporcionando pipelines preconstruidos para procesamiento de documentos.

## Frameworks de Análisis de Documentos

**deepdoctection** <sup>[13]</sup> es un framework de código abierto en Python que orquesta la extracción de documentos y tareas de análisis de layout usando modelos de aprendizaje profundo. No implementa sus propios modelos, pero permite construir pipelines que aprovechan bibliotecas altamente consideradas.

**PaddleOCR 3.0 Technical Report** <sup>[14]</sup> introduce un toolkit de código abierto bajo licencia Apache para OCR y análisis de documentos, presentando PP-OCRV5 para reconocimiento de texto multilingüe, PP-StructureV3 para análisis jerárquico de documentos, y PP-ChatOCRV4 para extracción de información clave.

## Mercado y Tendencias de la Industria

### Análisis de Mercado

**File Conversion Software Market** <sup>[15]</sup> estima que el mercado alcanzó USD 1.5 mil millones en 2024 y se proyecta que llegue a USD 3.1 mil millones para 2033, con un CAGR del 9.3%. El crecimiento está impulsado por la digitalización, necesidad de compatibilidad entre plataformas y adopción de modelos de trabajo remoto.

**Data Conversion Services Market** <sup>[16]</sup> proyecta crecimiento de USD 9.12 mil millones en 2024 a USD 14.56 mil millones con un CAGR del 5.48% durante 2025-2032, impulsado por la transformación digital y necesidades de migración de datos.

### Transformación Digital

**Digital Transformation Market** <sup>[17]</sup> muestra un crecimiento del mercado de USD 1.55 billones en 2024 a USD 15.82 billones estimados para 2034, con un CAGR del 26.15%. La transformación digital está impulsando la adopción de herramientas de conversión de documentos.

**Document Management Software Market** <sup>[18]</sup> se espera que crezca a un CAGR del 12.0% desde 2024 hasta 2034, alcanzando USD 24,322.8 millones, impulsado por la adopción de soluciones digitales y necesidades de gestión de documentos.

## Comparativas de Herramientas OCR

### Evaluación de Herramientas OCR de Código Abierto

**Which OCR toolset is good and why** <sup>[19]</sup> proporciona una comparación integral de herramientas OCR tanto propietarias como de código abierto. Los resultados muestran que herramientas propietarias como Nuance OmniPage pueden lograr 100% de precisión, mientras que en la categoría de código abierto, Tesseract (versiones 3 y 4.0) supera a GOCR y Ocrad.

**A Comparative Study of PDF Parsing Tools** <sup>[20]</sup> realiza una comparación exhaustiva de 10 herramientas de código abierto bien mantenidas para extracción de texto y tablas, utilizando el dataset DocLayNet. Los resultados muestran variaciones significativas en rendimiento según el tipo de documento.

### Avances en Formatos de Archivo

#### Formatos de Próxima Generación

**Next Gen Formats** <sup>[21]</sup> identifica WebP y AVIF como los dos formatos que actualmente califican como "próxima generación". WebP, creado por Google, pretende reemplazar JPEG, PNG y GIF, ofreciendo compresión con y sin pérdida y animación. AVIF utiliza AV1 junto con HEIF como contenedor.

**Future File Formats** <sup>[22]</sup> presenta un proyecto que busca desarrollar un formato de almacenamiento columnar de código abierto de próxima generación que se esfuerza por decodificación de alto rendimiento en hardware avanzado y alta portabilidad.

### Investigación en Procesamiento de Documentos

#### Modelos de Lenguaje Multimodal

**SmolDocling: An ultra-compact vision-language model** <sup>[23]</sup> introduce un modelo de visión-lenguaje ultra-compacto de 256M parámetros dirigido a conversión de documentos end-to-end. Genera DocTags, un nuevo formato de marcado universal que captura todos los elementos de página en su contexto completo con ubicación.

**DocGenome: An Open Large-scale Scientific Document Benchmark** <sup>[24]</sup> presenta un benchmark de documentos estructurados construido anotando 500K documentos científicos de 153 disciplinas de la comunidad de acceso abierto arXiv, proporcionando datos ground truth para entrenar y evaluar modelos multimodales.

### Tendencias Futuras y Direcciones de Investigación

## Inteligencia Artificial y Automatización

**Intelligent Document Processing Trends** <sup>[25]</sup> identifica siete tendencias clave para 2025, incluyendo OCR y NLP avanzados, IA agéntica, procesamiento de documentos adaptativo, y flujos de trabajo human-in-the-loop. La demanda de herramientas IDP seguras y compatibles está aumentando.

**E2E Process Automation Leveraging Generative AI** <sup>[26]</sup> examina conceptos y tendencias tecnológicas de IA generativa, OCR/IDP, RPA y agentes de automatización, analizando el potencial para automatización end-to-end de procesos de negocio.

## Herramientas de Conversión Emergentes

**Unstract Open Source Document Extraction** <sup>[27]</sup> presenta una herramienta de procesamiento de documentos de código abierto diseñada para extraer datos estructurados de documentos no estructurados, proporcionando un marco modular y flexible que permite personalización completa.

**DataDock: An Open Source Data Hub for Research** <sup>[28]</sup> aborda la escasez de servicios de investigación de código abierto de calidad con respecto a datos, proporcionando funciones especializadas para compartir, transferir y revisar archivos para equipos de investigación.

## Conclusiones

El panorama actual de la conversión de documentos y formatos de archivo está experimentando una revolución impulsada por la inteligencia artificial y tecnologías de código abierto. Las fuentes académicas más valoradas y recientes destacan:

1. **Docling** emerge como la herramienta más innovadora y completa para conversión de documentos con IA <sup>[1]</sup>
2. **Benchmarking riguroso** se está estableciendo para evaluar herramientas de procesamiento de documentos <sup>[2] [3]</sup>
3. **Formatos de próxima generación** como WebP, AVIF y NGFF están transformando el paisaje de formatos <sup>[6] [21]</sup>
4. **Mercado en crecimiento** con proyecciones de billones de dólares en los próximos años <sup>[15] [16] [17]</sup>
5. **Herramientas open source** cada vez más sofisticadas y competitivas <sup>[12] [10] [14]</sup>

La investigación académica continúa impulsando innovaciones en procesamiento inteligente de documentos, con énfasis en eficiencia, precisión y accesibilidad para la comunidad de código abierto.

✱

1. <https://arxiv.org/html/2501.17887v1>

2. <https://arxiv.org/html/2412.02210v3>

3. <https://arxiv.org/html/2412.07626v1>

4. <https://arxiv.org/html/2410.21169v2>
5. <https://arxiv.org/html/2412.19840v1>
6. <https://www.biorxiv.org/content/10.1101/2021.03.31.437929v4.full.pdf>
7. <https://nishtahir.com/the-wasm-text-format/>
8. <https://arxiv.org/pdf/1510.07851.pdf>
9. <https://pdf.abbyy.com/learning-center/pdf-standards/>
10. <https://arxiv.org/html/2310.10050>
11. <https://arxiv.org/abs/2310.01206>
12. <https://arxiv.org/html/2505.07672v1>
13. <https://konfuzio.com/en/deepdoctection/>
14. <https://arxiv.org/html/2507.05595v1>
15. <https://www.marketresearchintellect.com/product/global-file-converter-software-market-size-forecast/>
16. <https://www.marketresearchfuture.com/reports/data-conversion-service-market-35919>
17. <https://www.cervicornconsulting.com/digital-transformation-market>
18. <https://www.futuremarketinsights.com/reports/document-management-software-market>
19. <https://pdfs.semanticscholar.org/7c4d/264119ff0096d41ec90f7db6a9c12cc05afb.pdf>
20. <https://arxiv.org/pdf/2410.09871.pdf>
21. <https://cloudinary.com/glossary/next-gen-formats>
22. <https://db.cs.cmu.edu/projects/future-file-formats/>
23. <https://arxiv.org/html/2503.11576v1>
24. <https://arxiv.org/html/2406.11633v1>
25. <https://xtract.io/blog/intelligent-document-processing-trends/>
26. <https://arxiv.org/pdf/2505.20733.pdf>
27. <https://unstrat.com/blog/open-source-document-data-extraction-with-unstrat-deepseek/>
28. <https://arxiv.org/html/2406.16880v2>