

PDF

to PDF ▶

Contents

- Using PyPDF
- Using Unstructured
- Using PDFMiner
- Using PyMuPDF

This covers how to load pdfs into a document format that we can use downstream.

Using PyPDF

Load PDF using `pypdf` into array of documents, where each document contains the page content and metadata with `page` number.

```
from langchain.document_loaders import PyPDFLoader

loader = PyPDFLoader("example_data/layout-parser-paper.pdf")
pages = loader.load_and_split()
```

```
pages[0]
```

```
Document(page_content='LayoutParser : A Uni\x0ced Toolkit for Deep\nLearning Based Document Image Analysis\nZejiang Shen1( \x00), Ruochen Zhang2, Melissa Dell3, Benjamin Charles Germain\nLee4, Jacob Carlson3, and Weining Li5\n1Allen Institute for AI\nshannons@allenai.org\n2Brown University\nruochen zhang@brown.edu\n3Harvard University\nfmelissadell,jacob carlson g@fas.harvard.edu\n4University of Washington\nnbcgl@cs.washington.edu\n5University of Waterloo\nnw422li@uwaterloo.ca\nAbstract. Recent advances in document image analysis (DIA) have been\nprimarily driven by the application of neural networks. Ideally, research\noutcomes could be easily deployed in production and extended for further\ninvestigation. However, various factors like loosely organized codebases\nand sophisticated model con\x0cigurations complicate the easy reuse of
```

[Skip to main content](#)

model\ndevelopment in disciplines like natural language processing and computer\vision, none of them are optimized for challenges in the domain of DIA.\nThis represents a major gap in the existing toolkit, as DIA is central to\nacademic research across a wide range of disciplines in the social sciences\nand humanities. This paper introduces LayoutParser , an open-source\nlibrary for streamlining the usage of DL in DIA research and applica-\ntions. The core LayoutParser library comes with a set of simple and\nintuitive interfaces for applying and customizing DL models for layout de-\ntection, character recognition, and many other document processing tasks.\nTo promote extensibility, LayoutParser also incorporates a community\nplatform for sharing both pre-trained models and full document digiti-\nization pipelines. We demonstrate that LayoutParser is helpful for both\nlightweight and large-scale digitization pipelines in real-world use cases.\nThe library is publicly available at <https://layout-parser.github.io> .\nKeywords: Document Image Analysis ·Deep Learning ·Layout Analysis\n·Character Recognition ·Open Source library ·Toolkit.\n1 Introduction\nDeep Learning(DL)-based approaches are the state-of-the-art for a wide range of\ndocument image analysis (DIA) tasks including document image classi-\nfication [11,arXiv:2103.15348v2 [cs.CV] 21 Jun 2021', lookup_str='', metadata={'source': 'example_data/layout-parser-paper.pdf', 'page': '0'}, lookup_index=0)

An advantage of this approach is that documents can be retrieved with page numbers.

```
from langchain.vectorstores import FAISS
from langchain.embeddings.openai import OpenAIEmbeddings

faiss_index = FAISS.from_documents(pages, OpenAIEmbeddings())
docs = faiss_index.similarity_search("How will the community be engaged?", k=2)
for doc in docs:
    print(str(doc.metadata["page"]) + ":", doc.page_content)
```

9: 10 Z. Shen et al.

Fig. 4: Illustration of (a) the original historical Japanese document with layout detection results and (b) a recreated version of the document image that achieves much better character recognition recall. The reorganization algorithm rearranges the tokens based on the their detected bounding boxes given a maximum allowed height.

4LayoutParser Community Platform

Another focus of LayoutParser is promoting the reusability of layout detection models and full digitization pipelines. Similar to many existing deep learning libraries, LayoutParser comes with a community model hub for distributing layout models. End-users can upload their self-trained models to the model hub, and these models can be loaded into a similar interface as the currently available LayoutParser pre-trained models. For example, the model trained on the News Navigator dataset [17] has been incorporated in the model hub.

Beyond DL models, LayoutParser also promotes the sharing of entire document digitization pipelines. For example, sometimes the pipeline requires the combination of multiple DL models to achieve better accuracy. Currently, pipelines are mainly described in academic papers and implementations are often not pub-

[Skip to main content](#)

For each shared pipeline, it has a dedicated project page, with links to the source code, documentation, and an outline of the approaches. A discussion panel is provided for exchanging ideas. Combined with the core LayoutParser library, users can easily build reusable components based on the shared pipelines and apply them to solve their unique problems.

5 Use Cases

The core objective of LayoutParser is to make it easier to create both large-scale and light-weight document digitization pipelines. Large-scale document processing

Efficient Data Annotation Customized Model Training Model Customization DIA Model Hub DIA Pipeline Sharing Community Platform Layout Detection Models Document Images

The Core LayoutParser Library OCR Module Storage & Visualization Layout Data Structure

Fig. 1: The overall architecture of LayoutParser. For an input document image, the core LayoutParser library provides a set of off-the-shelf tools for layout detection, OCR, visualization, and storage, backed by a carefully designed layout data structure. LayoutParser also supports high level customization via efficient layout annotation and model training functions. These improve model accuracy on the target samples. The community platform enables the easy sharing of DIA models and whole digitization pipelines to promote reusability and reproducibility. A collection of detailed documentation, tutorials and exemplar projects make LayoutParser easy to learn and use.

AllenNLP [8] and transformers [34] have provided the community with complete DL-based support for developing and deploying models for general computer vision and natural language processing problems. LayoutParser, on the other hand, specializes specifically in DIA tasks. LayoutParser is also equipped with a community platform inspired by established model hubs such as Torch Hub [23] and TensorFlow Hub [1]. It enables the sharing of pretrained models as well as full document processing pipelines that are unique to DIA tasks.

There have been a variety of document data collections to facilitate the development of DL models. Some examples include PRIMA [3](magazine layouts), PubLayNet [38](academic paper layouts), Table Bank [18](tables in academic papers), Newspaper Navigator Dataset [16,17](newspaper figure layouts) and HJDataset [31](historical Japanese document layouts). A spectrum of models trained on these datasets are currently available in the LayoutParser model zoo to support different use cases.

3 The Core LayoutParser Library

At the core of LayoutParser is an off-the-shelf toolkit that streamlines DL-based document image analysis. Five components support a simple interface with comprehensive functionalities: 1) The layout detection models enable using pre-trained or self-trained DL models for layout detection with just four lines of code. 2) The detected layout information is stored in carefully engineered

Using Unstructured

```
from langchain.document_loaders import UnstructuredPDFLoader
```

[Skip to main content](#)

```
loader = UnstructuredPDFLoader("example_data/layout-parser-paper.pdf")
```

```
data = loader.load()
```

Retain Elements

Under the hood, Unstructured creates different “elements” for different chunks of text. By default we combine those together, but you can easily keep that separation by specifying

```
mode="elements".
```

```
loader = UnstructuredPDFLoader("example_data/layout-parser-paper.pdf",  
mode="elements")
```

```
data = loader.load()
```

```
data[0]
```

```
Document(page_content='LayoutParser: A Unified Toolkit for Deep\nLearning Based  
Document Image Analysis\nZejiang Shen1 (✉), Ruochen Zhang2, Melissa Dell3,  
Benjamin Charles Germain\nLee4, Jacob Carlson3, and Weining Li5\n1 Allen Institute  
for AI\nshannons@allenai.org\n2 Brown University\nruochen zhang@brown.edu\n3  
Harvard University\n{melissadell,jacob carlson}@fas.harvard.edu\n4 University of  
Washington\nbcgl@cs.washington.edu\n5 University of  
Waterloo\nnw422li@uwaterloo.ca\nAbstract. Recent advances in document image  
analysis (DIA) have been\nprimarily driven by the application of neural networks.  
Ideally, research\noutcomes could be easily deployed in production and extended  
for further\ninvestigation. However, various factors like loosely organized  
codebases\nand sophisticated model configurations complicate the easy reuse of im-  
\nportant innovations by a wide audience. Though there have been on-going\nefforts  
to improve reusability and simplify deep learning (DL) model\ndevelopment in  
disciplines like natural language processing and computer\nvision, none of them  
are optimized for challenges in the domain of DIA.\nThis represents a major gap in  
the existing toolkit, as DIA is central to\nacademic research across a wide range  
of disciplines in the social sciences\nand humanities. This paper introduces  
LayoutParser, an open-source\nlibrary for streamlining the usage of DL in DIA  
research and applica-\ntions. The core LayoutParser library comes with a set of  
simple and\nintuitive interfaces for applying and customizing DL models for layout  
de-\ntection, character recognition, and many other document processing tasks.\nTo  
promote extensibility, LayoutParser also incorporates a community\nplatform for
```

[Skip to main content](#)

```
digitization pipelines in real-world use cases.\n
The library is publicly available at https://layout-parser.github.io.\n
Keywords: Document Image Analysis · Deep Learning · Layout Analysis\n
· Character Recognition · Open Source library · Toolkit.\n
1\n
Introduction\n
Deep Learning(DL)-based approaches are the state-of-the-art for a wide range of\ndocument image analysis (DIA) tasks including document image classification [11,\narXiv:2103.15348v2 [cs.CV] 21 Jun 2021\n',
lookup_str='', metadata={'file_path': 'example_data/layout-parser-paper.pdf',
'page_number': 1, 'total_pages': 16, 'format': 'PDF 1.5', 'title': '', 'author': '',
'subject': '', 'keywords': '', 'creator': 'LaTeX with hyperref', 'producer': 'pdfTeX-1.40.21',
'creationDate': 'D:20210622012710Z', 'modDate': 'D:20210622012710Z', 'trapped': '',
'encryption': None}, lookup_index=0)
```

Fetching remote PDFs using Unstructured

This covers how to load online pdfs into a document format that we can use downstream. This can be used for various online pdf sites such as

<https://open.umn.edu/opentextbooks/textbooks/> and <https://arxiv.org/archive/>

Note: all other pdf loaders can also be used to fetch remote PDFs, but `OnlinePDFLoader` is a legacy function, and works specifically with `UnstructuredPDFLoader`.

```
from langchain.document_loaders import OnlinePDFLoader
```

```
loader = OnlinePDFLoader("https://arxiv.org/pdf/2302.03803.pdf")
```

```
data = loader.load()
```

```
print(data)
```

```
[Document(page_content='A WEAK ( k, k ) -LEFSCHETZ THEOREM FOR PROJECTIVE TORIC ORBIFOLDS\n\nWilliam D. Montoya\n\nInstituto de Matemática, Estatística e Computação Científica,\n\nIn [3] we proved that, under suitable conditions, on a very general codimension s quasi-smooth intersection subvariety X in a projective toric orbifold P d Σ with d + s = 2 ( k + 1 ) the Hodge conjecture holds, that is, every ( p, p ) -cohomology class, under the Poincaré duality is a rational linear combination of fundamental classes of algebraic subvarieties of X . The proof of the above-mentioned result relies, for p ≠ d + 1 - s , on a Lefschetz\n\nKeywords: (1,1)- Lefschetz theorem, Hodge conjecture, toric varieties, complete intersection\nEmail: wmontoya@ime.unicamp.br\n\ntheorem ([7]) and the Hard Lefschetz theorem for
```

[Skip to main content](#)

vector bundle, and the Cayley Proposition (4.3) which gives an isomorphism of some primitive cohomologies (4.2) of X and Y . The Cayley trick, following the philosophy of Mavlyutov in [7], reduces results known for quasi-smooth hypersurfaces to quasi-smooth intersection subvarieties. The idea in this paper goes the other way around, we translate some results for quasi-smooth intersection subvarieties to

Acknowledgement. I thank Prof. Ugo Bruzzo and Tiago Fonseca for useful discussions. I also acknowledge support from FAPESP postdoctoral grant No. 2019/23499-7.

Let M be a free abelian group of rank d , let $N = \text{Hom}(M, \mathbb{Z})$, and $N \otimes \mathbb{R} = N \otimes \mathbb{Z} \otimes \mathbb{R}$. If there exist k linearly independent primitive elements $e_1, \dots, e_k \in N$ such that $\sigma = \{ \mu_1 e_1 + \dots + \mu_k e_k \}$.

- The generators e_i are integral if for every i and any nonnegative rational number μ the product μe_i is in N only if μ is an integer.
- Given two rational simplicial cones σ, σ' one says that σ' is a face of σ ($\sigma' < \sigma$) if the set of integral generators of σ' is a subset of the set of integral generators of σ .
- A finite set $\Sigma = \{ \sigma_1, \dots, \sigma_t \}$ of rational simplicial cones is called a rational simplicial complete d -dimensional fan if:
 - all faces of cones in Σ are in Σ ;
 - if $\sigma, \sigma' \in \Sigma$ then $\sigma \cap \sigma' < \sigma$ and $\sigma \cap \sigma' < \sigma'$;
 - $N \otimes \mathbb{R} = \sigma_1 \cup \dots \cup \sigma_t$.

A rational simplicial complete d -dimensional fan Σ defines a d -dimensional toric variety P^d_Σ having only orbifold singularities which we assume to be projective. Moreover, $T := N \otimes \mathbb{Z} \otimes \mathbb{C}^* \simeq (\mathbb{C}^*)^d$ is the torus action on P^d_Σ . We denote by $\Sigma(i)$ the i -dimensional cones.

For a cone $\sigma \in \Sigma$, $\hat{\sigma}$ is the set of 1-dimensional cone in Σ that are not contained in σ and $x^{\hat{\sigma}} := \prod_{\rho \in \hat{\sigma}} x_\rho$ is the associated monomial in S .

Definition 2.2. The irrelevant ideal of P^d_Σ is the monomial ideal $B_\Sigma := \langle x^{\hat{\sigma}} \mid \sigma \in \Sigma \rangle$ and the zero locus $Z(\Sigma) := V(B_\Sigma)$ in the affine space $A^d := \text{Spec}(S)$ is the irrelevant locus.

Proposition 2.3 (Theorem 5.1.11 [5]). The toric variety P^d_Σ is a categorical quotient $A^d \setminus Z(\Sigma)$ by the group $\text{Hom}(\text{Cl}(\Sigma), \mathbb{C}^*)$ and the group action is induced by the $\text{Cl}(\Sigma)$ -grading of S .

Now we give a brief introduction to complex orbifolds and we mention the needed theorems for the next section. Namely: de Rham theorem and Dolbeault theorem for complex orbifolds.

Definition 2.4. A complex orbifold of complex dimension d is a singular complex space whose singularities are locally isomorphic to quotient singularities \mathbb{C}^d / G , for finite subgroups $G \subset \text{GL}(d, \mathbb{C})$.

Definition 2.5. A differential form on a complex orbifold Z is defined locally at $z \in Z$ as a G -invariant differential form on \mathbb{C}^d where $G \subset \text{GL}(d, \mathbb{C})$ and Z is locally isomorphic to \mathbb{C}^d .

Roughly speaking the local geometry of orbifolds reduces to local G -invariant geometry.

We have a complex of differential forms $(A^\bullet(Z), d)$ and a double complex $(A^\bullet, \bullet(Z), \partial, \bar{\partial})$ of bigraded differential forms which define the de Rham and the Dolbeault cohomology groups (for a fixed $p \in \mathbb{N}$) respectively:

(1,1)-Lefschetz theorem for projective toric orbifolds

Definition 3.1. A subvariety $X \subset P^d_\Sigma$ is quasi-smooth if $V(I_X) \subset A^d \setminus \Sigma(1)$ is smooth outside

Example 3.2. Quasi-smooth hypersurfaces or more generally quasi-smooth intersection sub-

Example 3.2. Quasi-smooth hypersurfaces or more generally quasi-smooth intersection sub-varieties are quasi-smooth subvarieties (see [2] or [7] for more details).

Remark 3.3. Quasi-smooth subvarieties are suborbifolds of P^d_Σ in the sense of Satake in [8]. Intuitively speaking they are subvarieties whose only singularities come from the ambient

Proof. From the exponential short exact sequence we have a long exact sequence in cohomology

$$H^1(0^*X) \rightarrow H^2(X, \mathbb{Z}) \rightarrow H^2(0^*X) \simeq H^0, 2(X)$$

where the last isomorphism is due to Steenbrink in [9]. Now, it is enough to prove the commutativity of the next diagram where the last isomorphism is due to Steenbrink in [9]. Now,

$$H^2(X, \mathbb{Z}) // H^2(X, 0^*X) \simeq \text{Dolbeault } H^2(X, \mathbb{C})^{\text{deRham}} \simeq H^2 dR(X, \mathbb{C}) // H^0, 2 \bar{\partial}(X)$$

the proof follows as the (1,1)-Lefschetz theorem in [6].

Remark 3.5. For $k =$

[Skip to main content](#)

[11] for details) we get an isomorphism of cohomologies :
 given by the Lefschetz morphism and since it is a morphism of Hodge structures, we have:
 $H^{1,1}(X, \mathbb{Q}) \cong H^{1,1}(X, \mathbb{Q})$
 Corollary 3.6. If the dimension of X is 1, 2 or 3. The Hodge conjecture holds on X .
 Proof. If the $\dim X = 1$ the result is clear by the Hard Lefschetz theorem for projective orbifolds. The dimension 2 and 3 cases are covered by Theorem 3.5 and the Hard Lefschetz.
 Cayley trick and Cayley proposition
 The Cayley trick is a way to associate to a quasi-smooth intersection subvariety a quasi-smooth hypersurface. Let L_1, \dots, L_s be line bundles on $P^d \Sigma$ and let $\pi : P(E) \rightarrow P^d \Sigma$ be the projective space bundle associated to the vector bundle $E = L_1 \oplus \dots \oplus L_s$. It is known that $P(E)$ is a $(d+s-1)$ -dimensional simplicial toric variety whose fan depends on the degrees of the line bundles and the fan Σ . Furthermore, if the Cox ring, without considering the grading, of $P^d \Sigma$ is $C[x_1, \dots, x_m]$ then the Cox ring of $P(E)$ is
 Moreover for X a quasi-smooth intersection subvariety cut off by f_1, \dots, f_s with $\deg(f_i) = [L_i]$ we relate the hypersurface Y cut off by $F = y_1 f_1 + \dots + y_s f_s$ which turns out to be quasi-smooth. For more details see Section 2 in [7].
 We will denote $P(E)$ as $P^{d+s-1} \Sigma, X$ to keep track of its relation with X and $P^d \Sigma$.
 The following is a key remark.
 Remark 4.1. There is a morphism $\iota : X \rightarrow Y \subset P^{d+s-1} \Sigma, X$. Moreover every point $z := (x, y) \in Y$ with $y \neq 0$ has a preimage. Hence for any subvariety $W = V(I_W) \subset X \subset P^d \Sigma$ there exists $W' \subset Y \subset P^{d+s-1} \Sigma, X$ such that $\pi(W') = W$, i.e., $W' = \{z = (x, y) \mid x \in W\}$.
 For $X \subset P^d \Sigma$ a quasi-smooth intersection variety the morphism in cohomology induced by the inclusion $i_* : H^{d-s}(P^d \Sigma, \mathbb{C}) \rightarrow H^{d-s}(X, \mathbb{C})$ is injective by Proposition 1.4 in [7].
 Definition 4.2. The primitive cohomology of $H^{d-s}_{\text{prim}}(X)$ is the quotient $H^{d-s}(X, \mathbb{C}) / i_*(H^{d-s}(P^d \Sigma, \mathbb{C}))$ and $H^{d-s}_{\text{prim}}(X, \mathbb{Q})$ with rational coefficients.
 $H^{d-s}(P^d \Sigma, \mathbb{C})$ and $H^{d-s}(X, \mathbb{C})$ have pure Hodge structures, and the morphism i_* is compatible with them, so that $H^{d-s}_{\text{prim}}(X)$ gets a pure Hodge structure.
 The next Proposition is the Cayley proposition.
 Proposition 4.3. [Proposition 2.3 in [3]] Let $X = X_1 \cap \dots \cap X_s$ be a quasi-smooth intersection subvariety in $P^d \Sigma$ cut off by homogeneous polynomials f_1, \dots, f_s . Then for $p \neq d+s-2, d+s-3$
 Remark 4.5. The above isomorphisms are also true with rational coefficients since $H^\bullet(X, \mathbb{C}) = H^\bullet(X, \mathbb{Q}) \otimes \mathbb{C}$. See the beginning of Section 7.1 in [10] for more details.
 Theorem 5.1. Let $Y = \{F = y_1 f_1 + \dots + y_k f_k = 0\} \subset P^{2k+1} \Sigma, X$ be the quasi-smooth hypersurface associated to the quasi-smooth intersection surface $X = X_{f_1} \cap \dots \cap X_{f_k} \subset P^{k+2} \Sigma$. Then on Y the Hodge conjecture holds.
 Proof. If $H^{k,k}_{\text{prim}}(X, \mathbb{Q}) = 0$ we are done. So let us assume $H^{k,k}_{\text{prim}}(X, \mathbb{Q}) \neq 0$. By the Cayley proposition $H^{k,k}_{\text{prim}}(Y, \mathbb{Q}) \cong H^{1,1}_{\text{prim}}(X, \mathbb{Q})$ and by the $(1,1)$ -Lefschetz theorem for projective toric orbifolds there is a non-zero algebraic basis $\lambda_{C_1}, \dots, \lambda_{C_n}$ with rational coefficients of $H^{1,1}_{\text{prim}}(X, \mathbb{Q})$, that is, there are $n := h^{1,1}_{\text{prim}}(X, \mathbb{Q})$ algebraic curves C_1, \dots, C_n in X such that under the Poincaré duality the class in homology $[C_i]$ goes to λ_{C_i} , $[C_i] \mapsto \lambda_{C_i}$. Recall that the Cox ring of P^{k+2} is contained in the Cox ring of $P^{2k+1} \Sigma, X$ without considering the grading. Considering the grading we have that if $\alpha \in Cl(P^{k+2} \Sigma)$ then $(\alpha, 0) \in Cl(P^{2k+1} \Sigma, X)$. So the polynomials defining $C_i \subset P^{k+2} \Sigma$ can be interpreted in $P^{2k+1} \Sigma, X$ but with different degree. Moreover, by Remark 4.1 each C_i is contained in $Y = \{F = y_1 f_1 + \dots + y_k f_k = 0\}$ and furthermore it has codimension k .
 Claim: $\{\lambda_{C_i}\}_{i=1}^n$ is a basis of $\text{prim}(\cdot)$. It is enough to prove that λ_{C_i} is different from zero in $H^{k,k}_{\text{prim}}(Y, \mathbb{Q})$ or equivalently that the cohomology classes λ_{C_i} do not come from the ambient space. By contradiction let

[Skip to main content](#)

($k + 2$) -dimensional algebraic subvariety $V \subset P^{2k+1} \Sigma, X$ such that $V \cap Y = C_j$ so they are equal as a homology class of $P^{2k+1} \Sigma, X$, i.e., $[V \cap Y] = [C_j]$. It is easy to check that $\pi(V) \cap X = C_j$ as a subvariety of $P^{k+2} \Sigma$ where $\pi: (x, y) \mapsto x$. Hence $[\pi(V) \cap X] = [C_j]$ which is equivalent to say that λC_j comes from $P^{k+2} \Sigma$ which contradicts the choice of $[C_j]$.

Remark 5.2. Into the proof of the previous theorem, the key fact was that on X the Hodge conjecture holds and we translate it to Y by contradiction. So, using an analogous argument we have:

Proposition 5.3. Let $Y = \{F = y_1 f_s + \dots + y_s f_s = 0\} \subset P^{2k+1} \Sigma, X$ be the quasi-smooth hypersurface associated to a quasi-smooth intersection subvariety $X = X_{f_1} \cap \dots \cap X_{f_s} \subset P^d \Sigma$ such that $d + s = 2(k + 1)$. If the Hodge conjecture holds on X then it holds as well on Y .

Corollary 5.4. If the dimension of Y is $2s - 1$, $2s$ or $2s + 1$ then the Hodge conjecture holds on Y .

Proof. By Proposition 5.3 and Corollary 3.6.

Angella, D. Cohomologies of certain orbifolds. Journal of Geometry and Physics

Batyrev, V. V., and Cox, D. A. On the Hodge structure of projective hypersurfaces in toric varieties. Duke Mathematical Journal

(Aug).

Bruzzo, U., and Montoya, W. On the Hodge conjecture for quasi-smooth intersections in toric varieties. S˜ao Paulo J. Math. Sci. Special Section: Geometry in Algebra and Algebra in Geometry

Caramello Jr, F. C. Introduction to orbifolds. arXiv:1105.3832

Cox, D., Little, J., and Schenck, H. Toric varieties, vol. American Mathematical Soc.

Griffiths, P., and Harris, J. Principles of Algebraic Geometry. John Wiley & Sons, Ltd.

Mavlyutov, A. R. Cohomology of complete intersections in toric varieties. Published in Pacific J. of Math.

No. 125 (1985).

Satake, I. On a Generalization of the Notion of Manifold. Proceedings of the National Academy of Sciences of the United States of America

(1957).

Steenbrink, J. H. M. Intersection form for quasi-homogeneous singularities. Compositio Mathematica

(1982).

Voisin, C. Hodge Theory and Complex Algebraic Geometry I, vol. of Cambridge Studies in Advanced Mathematics. Cambridge University Press.

Wang, Z. Z., and Zaffran, D. A remark on the Hard Lefschetz theorem for Kˆahler orbifolds. Proceedings of the American Mathematical Society

(Aug).

[2] Batyrev, V. V., and Cox, D. A. On the Hodge structure of projective hypersurfaces in toric varieties. Duke Mathematical Journal 75, 2 (Aug 1994).

[3] Bruzzo, U., and Montoya, W. On the Hodge conjecture for quasi-smooth intersections in toric varieties. S˜ao Paulo J. Math. Sci. Special Section: Geometry in Algebra and Algebra in Geometry

(2021).

A. R. Cohomology of complete intersections in toric varieties. Published online by Cambridge University Press, lookup_str='', metadata={'source': '/var/folders/ph/hhm7_zyx4l13k3v8z02dwp1w0000gn/T/tmpgq0ckaja/online_file.pdf'}, lookup_index=0]

Using PDFMiner

```
from langchain.document_loaders import PDFMinerLoader
```

[Skip to main content](#)


```
loader = PDFMinerLoader("example_data/layout-parser-paper.pdf")
```

```
data = loader.load()
```

Using PyMuPDF

This is the fastest of the PDF parsing options, and contains detailed metadata about the PDF and its pages, as well as returns one document per page.

```
from langchain.document_loaders import PyMuPDFLoader
```

```
loader = PyMuPDFLoader("example_data/layout-parser-paper.pdf")
```

```
data = loader.load()
```

```
data[0]
```

```
Document(page_content='LayoutParser: A Unified Toolkit for Deep\nLearning Based\nDocument Image Analysis\nZejiang Shen1 (✉), Ruochen Zhang2, Melissa Dell3,\nBenjamin Charles Germain\nLee4, Jacob Carlson3, and Weining Li5\n1 Allen Institute\nfor AI\nshannons@allenai.org\n2 Brown University\nruochen.zhang@brown.edu\n3\nHarvard University\n{melissadell,jacob.carlson}@fas.harvard.edu\n4 University of\nWashington\nnbcgl@cs.washington.edu\n5 University of\nWaterloo\nnw422li@uwaterloo.ca\nAbstract. Recent advances in document image\nanalysis (DIA) have been\nprimarily driven by the application of neural networks.\nIdeally, research\noutcomes could be easily deployed in production and extended\nfor further\ninvestigation. However, various factors like loosely organized\ncodebases\nand sophisticated model configurations complicate the easy reuse of im-\nportant innovations by a wide audience. Though there have been on-going\nefforts\nto improve reusability and simplify deep learning (DL) model\ndevelopment in\ndisciplines like natural language processing and computer\nvision, none of them\nare optimized for challenges in the domain of DIA.\nThis represents a major gap in\nthe existing toolkit, as DIA is central to\nacademic research across a wide range\nof disciplines in the social sciences\nand humanities. This paper introduces\nLayoutParser, an open-source\nlibrary for streamlining the usage of DL in DIA\nresearch and applica-\ntions. The core LayoutParser library comes with a set of\nsimple and\nintuitive interfaces for applying and customizing DL models for layout\nde-\ntection, character recognition, and many other document processing tasks.\nTo
```

[Skip to main content](#)

```
demonstrate that LayoutParser is helpful for both\n\nlightweight and large-scale
digitization pipelines in real-world use cases.\n\nThe library is publicly available
at https://layout-parser.github.io.\n\nKeywords: Document Image Analysis · Deep
Learning · Layout Analysis\n\n· Character Recognition · Open Source library ·
Toolkit.\n\n1\n\nIntroduction\n\nDeep Learning(DL)-based approaches are the state-of-the-
art for a wide range of\ndocument image analysis (DIA) tasks including document
image classification [11,\narXiv:2103.15348v2 [cs.CV] 21 Jun 2021\n',
lookup_str='', metadata={'file_path': 'example_data/layout-parser-paper.pdf',
'page_number': 1, 'total_pages': 16, 'format': 'PDF 1.5', 'title': '', 'author':
'', 'subject': '', 'keywords': '', 'creator': 'LaTeX with hyperref', 'producer':
'pdfTeX-1.40.21', 'creationDate': 'D:20210622012710Z', 'modDate':
'D:20210622012710Z', 'trapped': '', 'encryption': None}, lookup_index=0)
```

Additionally, you can pass along any of the options from the [PyMuPDF documentation](#) as keyword arguments in the `load` call, and it will be pass along to the `get_text()` call.