

INDEPENDENT RESEARCH IN DATA SCIENCE

Entity Extraction Strategies for Historical Mortgage Records

By Uday Sapra

Under Kate Thomas and Jacob Faber (NYU Redlining Lab)

Dec 12, 2024

Agenda

- What was the research **about**?
- What was my **role**?
- What has been the **output**?
- What has been **learnt**?
- What is the **approach** going forward?

01 The Team & Research

02 Literature Review

03 Incompatibility with our Data

04 Renewed Approach

05 Outputs and Validation

06 Conclusion and Next Steps

The Research

Data-Oriented Approach to Delineating the Long Term Impacts of Government Redlining

- **Redlining in Real Estate:** Denying access to credit (mortgage and refinancing) based on race
- **Impact:** Segregation, inequalities in neighbourhood development, wealth creation, opportunities, & asset values

My Role?

Create the **Process** to build the Dataset that enables this research - from **Raw Document Scans**

OUR DATA

280,000

Mortgage Document Scans (.TIF)

20

Diverse Counties

1930-1975

Four Decades Spanned

Document Corpus

FHA Form No. 2151-b
(For use after Section 802)
(Revised Mar. 1, 1968)

DEED OF TRUST

STATE OF TEXAS
COUNTY OF GALVESTON

THIS INSTRUMENT, made and entered into by and between L. B. Berndt and wife, Mary Ellen Berndt, of the County of Galveston in the State of Texas, hereinafter called the Grantors, and Wm. S. Bradley, Trustee(s), of Dallas, Texas, hereinafter called the Trustee

WITNESSETH That the Grantors for and in consideration of the sum of Ten Dollars (\$10.00) and other valuable consideration in hand paid, the receipt whereof is hereby acknowledged, and the further consideration, uses, purposes and trusts herein set forth and declared, have granted, bargained, sold and conveyed, and by these presents do grant, bargain, sell and convey unto the said Trustee, and unto his successors in the trust hereby created and his assigns, forever, all of the following described real estate together with all the improvements thereon and hereafter placed thereon situated in the County of Galveston, State of Texas, to-wit:

Lot Nineteen (19) of Palm Gardens, in the City and County of Galveston, Texas, according to map of said Palm Gardens of record in Volume 417, page 265, in the office of the County Clerk of Galveston County, Texas.

To HAVE AND TO HOLD the above-described premises, together with all the rights, hereditaments, and appurtenances in anywise appertaining or belonging thereto, including all heating, plumbing, refrigeration and lighting fixtures and equipment now or hereafter attached thereto or used in connection therewith, unto the said Trustee, his successors in this trust and his assigns, forever. And the Grantors do hereby bind themselves and their heirs, executors, administrators, and legal representatives, to warrant and forever defend all and singular the said premises unto the said Trustee, and unto his successors in this trust, and his assigns, forever, against any person who lawfully claims or shall claim the same or any part thereof

This conveyance is made in trust to secure the payment of the principal sum of FORTY-SIX HUNDRED AND NO/100 Dollars (\$460.00), as evidenced by a certain promissory note of even date herewith executed by the Grantors, payable to the order of MORTGAGE INVESTMENT CORPORATION, the terms of which are incorporated herein by reference, together with interest at the rate of four and one-half per centum (4 1/2%) per annum on the unpaid balance, both interest and principal being payable monthly as it accrues at the office of Mortgage Investment Corporation in Dallas, Texas, in monthly installments of Twenty-five and 58/100 Dollars (\$25.58) each, including interest, one on the first day of each month hereafter, commencing on the first day of April, 1940, and continuing until the principal and interest are fully paid, except that the final payment of principal and interest, if not sooner paid, shall be due and payable on the first day of March, 1965. The note also provides that if default is made in the payment of any installment thereunder, and if the default is not made good prior to the due date of the next such installment, at the option of the holder, the note shall become immediately due and payable without notice, and that this lien may be foreclosed. Failure to exercise this option is not to constitute a waiver of the right to exercise it in the event of any subsequent default. If the note is placed in the hands of an attorney for collection, or is collected through the Probate Court or the Bankruptcy Court or through other legal proceedings, the makers thereof agree to pay, as attorney's fees, an additional amount equal to ten per centum (10%) of the amount then owing on the note

The Grantors covenant as follows:

1. That they will pay the principal of and interest on the note secured hereby in accordance with the terms thereof. Privilege is reserved to pay the debt in whole, or in an amount equal to one or more monthly payments on the principal that are next due on the note, on the first day of any month prior to maturity; provided, however, that written notice of an intention to exercise such privilege is given at least thirty (30) days prior to prepayment; and provided further that in the event the debt is paid in full prior to maturity and at that time it is insured under the provisions of the National Housing Act, they will pay to the holder of the note an adjusted premium charge of one per centum (1%) of the original principal amount thereof, except that in no event shall the adjusted premium exceed the aggregate amount of premium charges which would have been payable if this Deed of Trust and the note secured hereby had continued to be insured until maturity, such payment to be applied by the holder of the note upon its obligation to the Federal Housing Administrator on account of mortgage insurance.

5. Grantors will keep the improvements now existing or hereafter erected on the said premises, insured as may be required from time to time by the holder of the note against loss by fire and other hazards, casualty and contingencies in such amounts and for such periods as may be required by said holder, and will pay promptly when due any premiums on such insurance for payment of which provision has not been made hereunder. All insurance shall be carried in companies approved by the holder and the policies and receipts therefor shall be held by said holder and have attached thereto non payable clauses in favor of and in form acceptable to the said holder.

6. Grantors will give immediate notice by mail to the holder of the note, and said holder may make proof of loss if not made promptly by Grantors, and each insurance company concerned is hereby authorized and directed to make payment for such loss directly to the holder instead of to Grantors and said holder, jointly, and the insurance proceeds, or any part thereof, shall be applied by the holder to the reduction of the indebtedness hereby secured or to the restoration or repair of the property damaged, in event of foreclosure of this deed of trust or other transfer of title to the said premises to satisfaction of the indebtedness secured hereby.

7. Grantors will permit holder of said note, its agents or representatives, to inspect the said premises at any time; will maintain the said premises free from waste or substance of any kind and in good condition, and make all repairs, replacements, improvements and additions which may be necessary to preserve and maintain the premises and the value thereof; will comply with all laws, ordinances and regulations and all covenants, conditions and restrictions of any kind, affecting said property or its use; will not erect, remove or alter any building, structure or other property covered by this deed of trust, or permit the same to be altered, constructed, removed or used for any purpose other than that for which it is now used, without first obtaining the permission in writing of the holder of the note; will complete in a good and workmanlike manner any building which is begun or may be commenced or repaired thereon, will pay when due all claims for labor performed and material furnished, and will not permit any lien of mechanics or materialmen to attach to said premises.

COUNTY OF GALVESTON /

T H I S I N D E N T U R E, made and entered into by and between KERMIT E. AGEE & MARGARET AGEE (Husband and Wife) of the County of Galveston in the State of Texas, hereinafter called the Grantors, and R. M. ORTH, Trustee, of Texas City, Texas, hereinafter called the Trustee:

WITNESSETH: That the Grantors, for and in consideration of the sum of Ten

1940, and continuing until the principal and interest are fully paid, except that the final payment of principal and interest, if not sooner paid, shall be due and payable on the first day of AUGUST, 1965. The note also provides that if

ST

HAYNES, JR. and ELEANOR HAYNES (Husband and Wife) of the State of Texas, hereinafter called the Grantors, and Ernest F. Wenzel, Trustee, of Dallas, Texas, hereinafter called the Trustee

in Dollars (\$10.00) and other valuable consideration in hand paid, the receipt whereof is hereby acknowledged, and the further consideration, uses, purposes and trusts herein set forth and declared, have granted, bargained, sold and conveyed, and by these presents do grant, bargain, sell and convey unto the said Trustee, and unto his successors in the trust hereby created and his assigns, forever, all of the following described real estate together with all the improvements thereon and hereafter placed thereon situated in the County of Galveston, State of Texas, to-wit:

Lot Nineteen (19) of Palm Gardens, in the City and County of Galveston, Texas, according to map of said Palm Gardens of record in Volume 417, page 265, in the office of the County Clerk of Galveston County, Texas.

To HAVE AND TO HOLD the above-described premises, together with all the rights, hereditaments, and appurtenances in anywise appertaining or belonging thereto, including all heating, plumbing, refrigeration and lighting fixtures and equipment now or hereafter attached thereto or used in connection therewith, unto the said Trustee, his successors in this trust and his assigns, forever. And the Grantors do hereby bind themselves and their heirs, executors, administrators, and legal representatives, to warrant and forever defend all and singular the said premises unto the said Trustee, and unto his successors in this trust, and his assigns, forever, against any person who lawfully claims or shall claim the same or any part thereof

This conveyance is made in trust to secure the payment of the principal sum of FORTY-SIX HUNDRED AND NO/100 Dollars (\$460.00), as evidenced by a certain promissory note of even date herewith executed by the Grantors, payable to the order of MORTGAGE INVESTMENT CORPORATION, the terms of which are incorporated herein by reference, together with interest at the rate of four and one-half per centum (4 1/2%) per annum on the unpaid balance, both interest and principal being payable monthly as it accrues at the office of Mortgage Investment Corporation in Dallas, Texas, in monthly installments of Twenty-five and 58/100 Dollars (\$25.58) each, including interest, one on the first day of each month hereafter, commencing on the first day of April, 1940, and continuing until the principal and interest are fully paid, except that the final payment of principal and interest, if not sooner paid, shall be due and payable on the first day of March, 1965. The note also provides that if default is made in the payment of any installment thereunder, and if the default is not made good prior to the due date of the next such installment, at the option of the holder, the note shall become immediately due and payable without notice, and that this lien may be foreclosed. Failure to exercise this option is not to constitute a waiver of the right to exercise it in the event of any subsequent default. If the note is placed in the hands of an attorney for collection, or is collected through the Probate Court or the Bankruptcy Court or through other legal proceedings, the makers thereof agree to pay, as attorney's fees, an additional amount equal to ten per centum (10%) of the amount then owing on the note

The Grantors covenant as follows:

1. That they will pay the principal of and interest on the note secured hereby in accordance with the terms thereof. Privilege is reserved to pay the debt in whole, or in an amount equal to one or more monthly payments on the principal that are next due on the note, on the first day of any month prior to maturity; provided, however, that written notice of an intention to exercise such privilege is given at least thirty (30) days prior to prepayment; and provided further that in the event the debt is paid in full prior to maturity and at that time it is insured under the provisions of the National Housing Act, they will pay to the holder of the note an adjusted premium charge of one per centum (1%) of the original principal amount thereof, except that in no event shall the adjusted premium exceed the aggregate amount of premium charges which would have been payable if this Deed of Trust and the note secured hereby had continued to be insured until maturity, such payment to be applied by the holder of the note upon its obligation to the Federal Housing Administrator on account of mortgage insurance.

5. Grantors will keep the improvements now existing or hereafter erected on the said premises, insured as may be required from time to time by the holder of the note against loss by fire and other hazards, casualty and contingencies in such amounts and for such periods as may be required by said holder, and will pay promptly when due any premiums on such insurance for payment of which provision has not been made hereunder. All insurance shall be carried in companies approved by the holder and the policies and receipts therefor shall be held by said holder and have attached thereto non payable clauses in favor of and in form acceptable to the said holder.

6. Grantors will give immediate notice by mail to the holder of the note, and said holder may make proof of loss if not made promptly by Grantors, and each insurance company concerned is hereby authorized and directed to make payment for such loss directly to the holder instead of to Grantors and said holder, jointly, and the insurance proceeds, or any part thereof, shall be applied by the holder to the reduction of the indebtedness hereby secured or to the restoration or repair of the property damaged, in event of foreclosure of this deed of trust or other transfer of title to the said premises to satisfaction of the indebtedness secured hereby.

7. Grantors will permit holder of said note, its agents or representatives, to inspect the said premises at any time; will maintain the said premises free from waste or substance of any kind and in good condition, and make all repairs, replacements, improvements and additions which may be necessary to preserve and maintain the premises and the value thereof; will comply with all laws, ordinances and regulations and all covenants, conditions and restrictions of any kind, affecting said property or its use; will not erect, remove or alter any building, structure or other property covered by this deed of trust, or permit the same to be altered, constructed, removed or used for any purpose other than that for which it is now used, without first obtaining the permission in writing of the holder of the note; will complete in a good and workmanlike manner any building which is begun or may be commenced or repaired thereon, will pay when due all claims for labor performed and material furnished, and will not permit any lien of mechanics or materialmen to attach to said premises.

COUNTY OF GALVESTON /

T H I S I N D E N T U R E, made and entered into by and between KERMIT E. AGEE & MARGARET AGEE (Husband and Wife) of the County of Galveston in the State of Texas, hereinafter called the Grantors, and R. M. ORTH, Trustee, of Texas City, Texas, hereinafter called the Trustee:

WITNESSETH: That the Grantors, for and in consideration of the sum of Ten

1940, and continuing until the principal and interest are fully paid, except that the final payment of principal and interest, if not sooner paid, shall be due and payable on the first day of AUGUST, 1965. The note also provides that if

App. No. 15777

between A. M. MOORE, JR. of the State of Texas, hereinafter called the Grantor, and Ernest F. Wenzel, Trustee, of Dallas County, Texas, hereinafter called the Trustee

in Dollars (\$10.00) and other valuable consideration in hand paid, the receipt whereof is hereby acknowledged, and the further consideration, uses, purposes and trusts herein set forth and declared, have granted, bargained, sold and conveyed, and by these presents do grant, bargain, sell and convey unto the said Trustee, and unto his successors in the trust hereby created and his assigns, forever, all of the following described real estate together with all the improvements thereon and hereafter placed thereon situated in the County of Galveston, State of Texas, to-wit:

Lot Nineteen (19) of Palm Gardens, in the City and County of Galveston, Texas, according to map of said Palm Gardens of record in Volume 417, page 265, in the office of the County Clerk of Galveston County, Texas.

To HAVE AND TO HOLD the above-described premises, together with all the rights, hereditaments, and appurtenances in anywise appertaining or belonging thereto, including all heating, plumbing, refrigeration and lighting fixtures and equipment now or hereafter attached thereto or used in connection therewith, unto the said Trustee, his successors in this trust and his assigns, forever. And the Grantors do hereby bind themselves and their heirs, executors, administrators, and legal representatives, to warrant and forever defend all and singular the said premises unto the said Trustee, and unto his successors in this trust, and his assigns, forever, against any person who lawfully claims or shall claim the same or any part thereof

This conveyance is made in trust to secure the payment of the principal sum of FORTY-SIX HUNDRED AND NO/100 Dollars (\$460.00), as evidenced by a certain promissory note of even date herewith executed by the Grantors, payable to the order of MORTGAGE INVESTMENT CORPORATION, the terms of which are incorporated herein by reference, together with interest at the rate of four and one-half per centum (4 1/2%) per annum on the unpaid balance, both interest and principal being payable monthly as it accrues at the office of Mortgage Investment Corporation in Dallas, Texas, in monthly installments of Twenty-five and 58/100 Dollars (\$25.58) each, including interest, one on the first day of each month hereafter, commencing on the first day of April, 1940, and continuing until the principal and interest are fully paid, except that the final payment of principal and interest, if not sooner paid, shall be due and payable on the first day of March, 1965. The note also provides that if default is made in the payment of any installment thereunder, and if the default is not made good prior to the due date of the next such installment, at the option of the holder, the note shall become immediately due and payable without notice, and that this lien may be foreclosed. Failure to exercise this option is not to constitute a waiver of the right to exercise it in the event of any subsequent default. If the note is placed in the hands of an attorney for collection, or is collected through the Probate Court or the Bankruptcy Court or through other legal proceedings, the makers thereof agree to pay, as attorney's fees, an additional amount equal to ten per centum (10%) of the amount then owing on the note

The Grantors covenant as follows:

1. That they will pay the principal of and interest on the note secured hereby in accordance with the terms thereof. Privilege is reserved to pay the debt in whole, or in an amount equal to one or more monthly payments on the principal that are next due on the note, on the first day of any month prior to maturity; provided, however, that written notice of an intention to exercise such privilege is given at least thirty (30) days prior to prepayment; and provided further that in the event the debt is paid in full prior to maturity and at that time it is insured under the provisions of the National Housing Act, they will pay to the holder of the note an adjusted premium charge of one per centum (1%) of the original principal amount thereof, except that in no event shall the adjusted premium exceed the aggregate amount of premium charges which would have been payable if this Deed of Trust and the note secured hereby had continued to be insured until maturity, such payment to be applied by the holder of the note upon its obligation to the Federal Housing Administrator on account of mortgage insurance.

5. Grantors will keep the improvements now existing or hereafter erected on the said premises, insured as may be required from time to time by the holder of the note against loss by fire and other hazards, casualty and contingencies in such amounts and for such periods as may be required by said holder, and will pay promptly when due any premiums on such insurance for payment of which provision has not been made hereunder. All insurance shall be carried in companies approved by the holder and the policies and receipts therefor shall be held by said holder and have attached thereto non payable clauses in favor of and in form acceptable to the said holder.

6. Grantors will give immediate notice by mail to the holder of the note, and said holder may make proof of loss if not made promptly by Grantors, and each insurance company concerned is hereby authorized and directed to make payment for such loss directly to the holder instead of to Grantors and said holder, jointly, and the insurance proceeds, or any part thereof, shall be applied by the holder to the reduction of the indebtedness hereby secured or to the restoration or repair of the property damaged, in event of foreclosure of this deed of trust or other transfer of title to the said premises to satisfaction of the indebtedness secured hereby.

7. Grantors will permit holder of said note, its agents or representatives, to inspect the said premises at any time; will maintain the said premises free from waste or substance of any kind and in good condition, and make all repairs, replacements, improvements and additions which may be necessary to preserve and maintain the premises and the value thereof; will comply with all laws, ordinances and regulations and all covenants, conditions and restrictions of any kind, affecting said property or its use; will not erect, remove or alter any building, structure or other property covered by this deed of trust, or permit the same to be altered, constructed, removed or used for any purpose other than that for which it is now used, without first obtaining the permission in writing of the holder of the note; will complete in a good and workmanlike manner any building which is begun or may be commenced or repaired thereon, will pay when due all claims for labor performed and material furnished, and will not permit any lien of mechanics or materialmen to attach to said premises.

COUNTY OF GALVESTON /

T H I S I N D E N T U R E, made and entered into by and between KERMIT E. AGEE & MARGARET AGEE (Husband and Wife) of the County of Galveston in the State of Texas, hereinafter called the Grantors, and R. M. ORTH, Trustee, of Texas City, Texas, hereinafter called the Trustee:

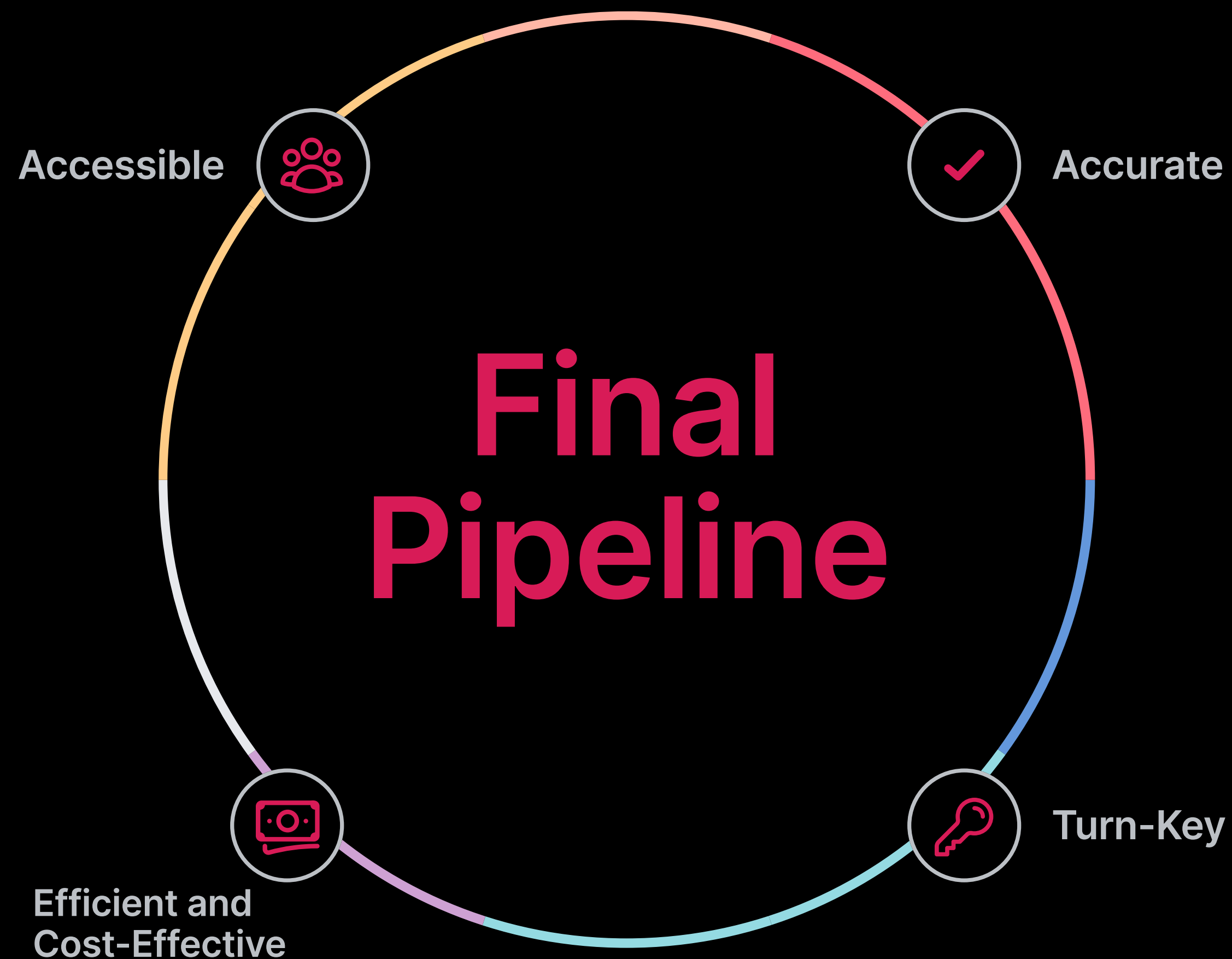
WITNESSETH: That the Grantors, for and in consideration of the sum of Ten

1940, and continuing until the principal and interest are fully paid, except that the final payment of principal and interest, if not sooner paid, shall be due and payable on the first day of AUGUST, 1965. The note also provides that if

CSV with 12 Features Per Doc

E	F	G	H	I	J
Second_Borrower	Other_Party_First	Other_Party_Last	Lending_Bank	Interest_Rate	Loan_Amount
N/A	J. G.	HESTWOOD	FIRST MORTGAGE COMPANY OF HOUSTON	4.50%	\$3900.00
N/A	J.G.	ESTWOOD	FIRST MORTGAGE COMPANY OF HOUSTON	4.50%	\$3550.00
Berndt	Wm. S.	Bradley	Mortgage Investment Corporation	4.50%	\$4600.00
N/A	R. M.	Orth	The Texas City National Bank	4.50%	\$3,000.00
Quinn	B. B.	Yeager	Gulf Coast Investment Corporation	4.50%	\$4,400.00
N/A	J.E.	Aosashe Jr.	J. E. FOSTER & SON, INC.	4.50%	\$3,950.00
N/A	J.	Hestwood	First Mortgage Company of Houston, Inc.	4.50%	\$4,950.00

Eventually the **First Large Scale** Public Dataset for FHA/VA Lending Research



The Objective

Develop an accurate yet **efficient, cost-effective,** and **turn-key** approach for digitizing real-estate documents **at scale.**

- Current Manual Labelling Budget: **\$300,000**
- Human Labeling Time Per Document: **5-10 Minutes**

Existing Literature

#1 **YOLO APPROACH** Manually Enter Data

#2 **COMPUTER VISION APPROACH** Document Layout Analysis + OCR + Entity Recognition

CV PIPELINE OVERVIEW

- ↓ Use CV Models to Identify Document Layout and Extract Areas of Interest
- ↓ Run OCR on Extracted Segments to get Text with Higher Accuracy
- ↓ Use SpaCY or LLMs for Extracting Entities from Segment Strings

3

Computer Vision Models Tested from Meta's Detectron Library

2

OCR Models Tested

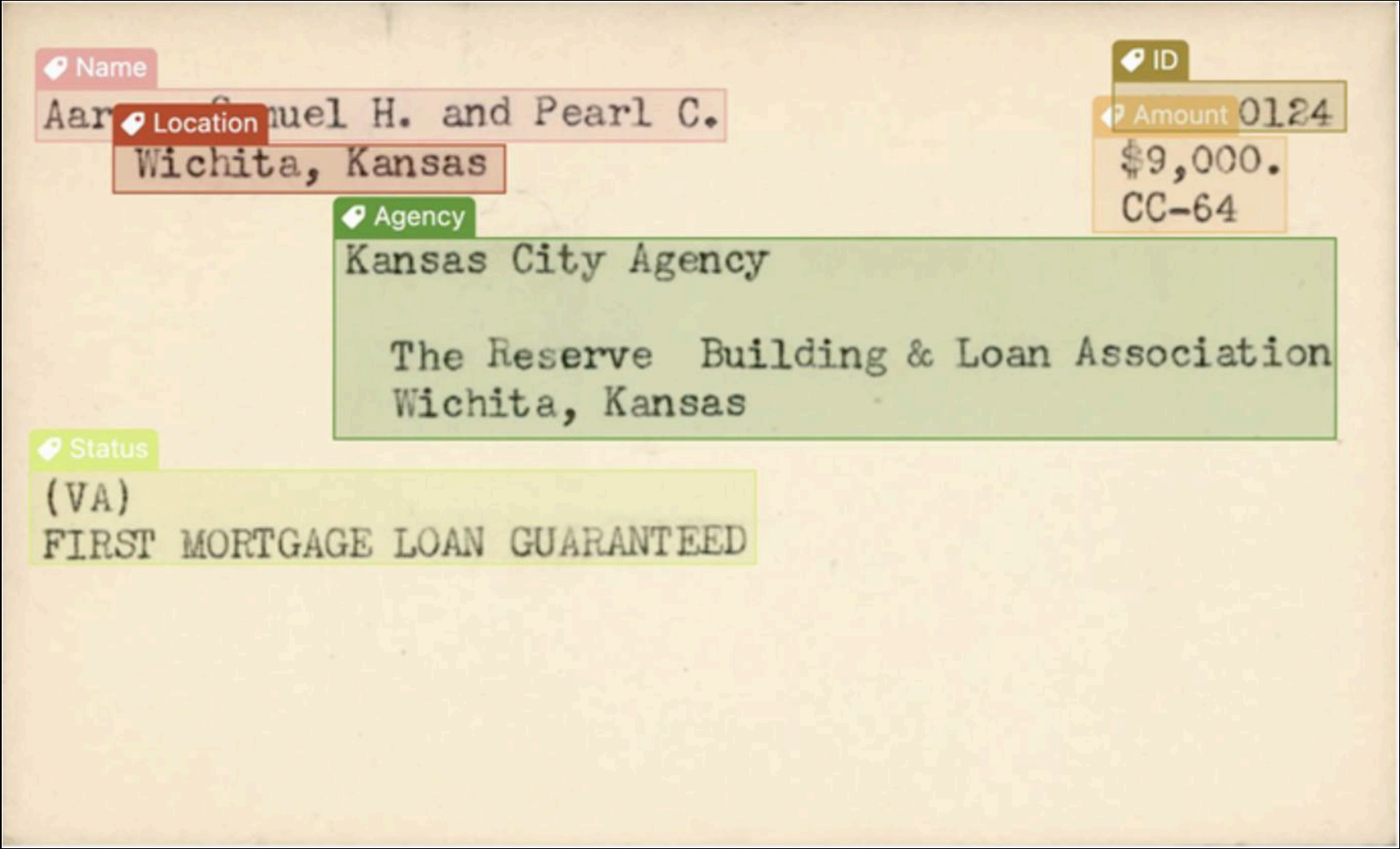
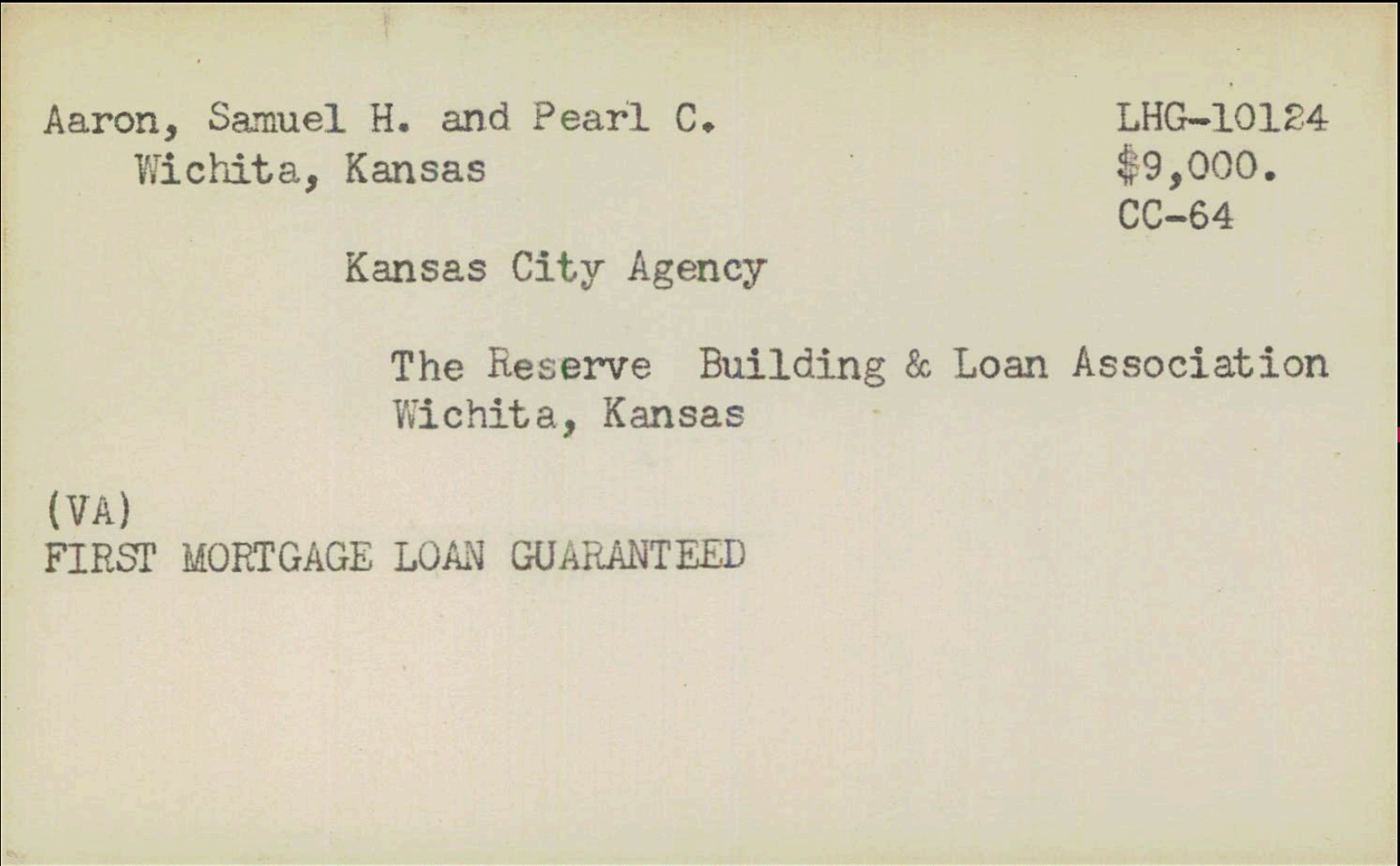
3

Entity Recognition Models Tested including HuggingFace models and SpaCy

Why did they do this?

- OCR models (Tesseract/Adobe) were not robust enough to parse entire documents.
- Entity Recognition models were not powerful enough entirely handle downstream processing.

Detectron models output **bounding boxes** and **entity label predictions**. This solved both issues - segmented the document for **better OCR accuracy** and classified boxes with entity labels to **reduce downstream processing**.



```
{
  "Text": "Kansas City Agency \n The Reserve Building & Loan Association Wichita, Kansas",
  "Label": "Agency"
}
```



Below are document samples from Galveston county for the year 1940. As can be seen, **these documents exhibit high layout variability and textual density.**

213

Said note being secured by vendor's lien retained in deed dated December 15, 1939, from Johnnie Chuoke to Ernest F. Wenzel and wife, Rosie Wenzel, of record in the office of the County Clerk of Galveston County, Texas, in Vol. 597, pages 574-75, and additionally secured by deed of trust dated December 15, 1939, from Ernest F. Wenzel and wife, Rosie Wenzel, to Maco Stewart, Trustee for Johnnie Chuoke, of record in said county clerk's office in Vol. 600, pages 329-30, on the following described real estate situate in the City & County of Galveston, Texas:

Lot Three (3) in Block One Hundred Forty-six (146) of Denver Resurvey, according to map of said Denver Resurvey of record in the office of the County Clerk of Galveston County, Texas, in Vol. 91, page 196,

in the City of Galveston, County of Galveston Texas.

And I do further Grant, Sell and Convey unto the said Thos. F. Davis all the rights, title, interest and liens owned or held by me in said land by virtue of said note herein conveyed and assigned.

TO HAVE AND TO HOLD unto the said THOS. F. DAVIS, his heirs and assigns, the above described note, together with all and singular the above mentioned liens and any and all liens, rights, equities, remedies, privileges, titles and interest in and to said land, which I have by virtue of being the legal holder and owner of said note; HEREBY WARRANTING that all payments, credits and offsets to which said note is entitled appear thereon, but without recourse on me. warranting, however, that there is due and unpaid the sum of \$2000.00.

EXECUTED this the 16th day of September, 1940.

Johnnie Chuoke

THE STATE OF TEXAS }
COUNTY OF GALVESTON }

BEFORE ME, the undersigned, a Notary Public in and for the State and County aforesaid, on this day personally appeared JOHNNIE CHUOKE, known to me to be the person whose name is subscribed to the foregoing instrument, and acknowledged to me that he executed the same for the purposes and consideration therein expressed.

Given Under My Hand and Seal of Office this the 16th day of September, 1940.

Henry Clark
(SEAL) Notary Public for Galveston County, Texas.

Filed for record September 16th., 1940, at 4:45 o'clock P. M.
Recorded September 17th., 1940, at 4:10 o'clock P. M.
John R. Platte, Clerk. By Edith H. Brown Deputy.

Kernit E.
Agee and wife,
To
R. M. Orth,
Trustee,
DEED OF TRUST

FHA - For Use Under Title I
Class 3 Loans on Leasehold Estates.

DEED OF TRUST

STATE OF TEXAS }
COUNTY OF GALVESTON } SS:

T H I S I N D E N T U R E, made and entered into by and between KERMIT E. AGEE & MARGARET AGEE (Husband and Wife) of the County of Galveston in the State of Texas, hereinafter called the Grantors, and R. M. ORTH, Trustee, of Texas City, Texas, hereinafter called the Trustee:

WITNESSETH: That the Grantors, for and in consideration of the sum of Ten

Renewed Approach



1

Modern OCR Models

Modern open and closed-source CNN & RNN based OCR models exist that excel at **word** recognition.

2

Parse Entire Documents

Parse entire documents **eliminating** the CV step for accuracy.

PyTorch Models, therefore GPU Acceleration.

3

Introduce LLMs

LLMs understand context.

They can also smooth and preprocess OCR errors.

Prompt-Engineering will help us tune LLMs for our extraction task.

4

Validate Output and Fine-Tune

Test different OCR models and LLMs.

Validate with relevant metrics.

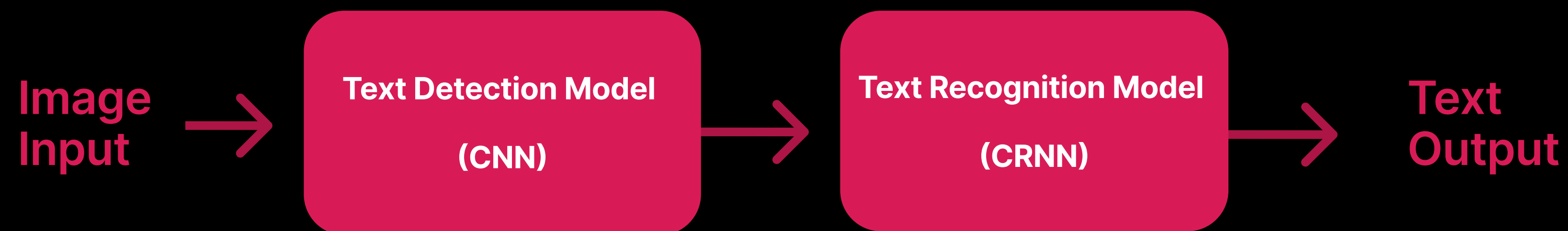
Analyze errors to find ideal combinations.



OCR Stage

LEVERAGING MODERN CRNN-BASED OCR MODELS

2-Step OCR Workflow



Library Choice

docTR

- Higher Accuracy
- Model Choices
- Better API

Model Choice

Text-Detection Model: **DB Resnet 50**

Text-Recognition Model: **CRNN VGG 16**
(VGG 16+Bi-Directional RNN)

Closest to Amazon Textract

Inference Speed: 19 Seconds per Doc (L4 GPU)

OCR MODELS TESTED

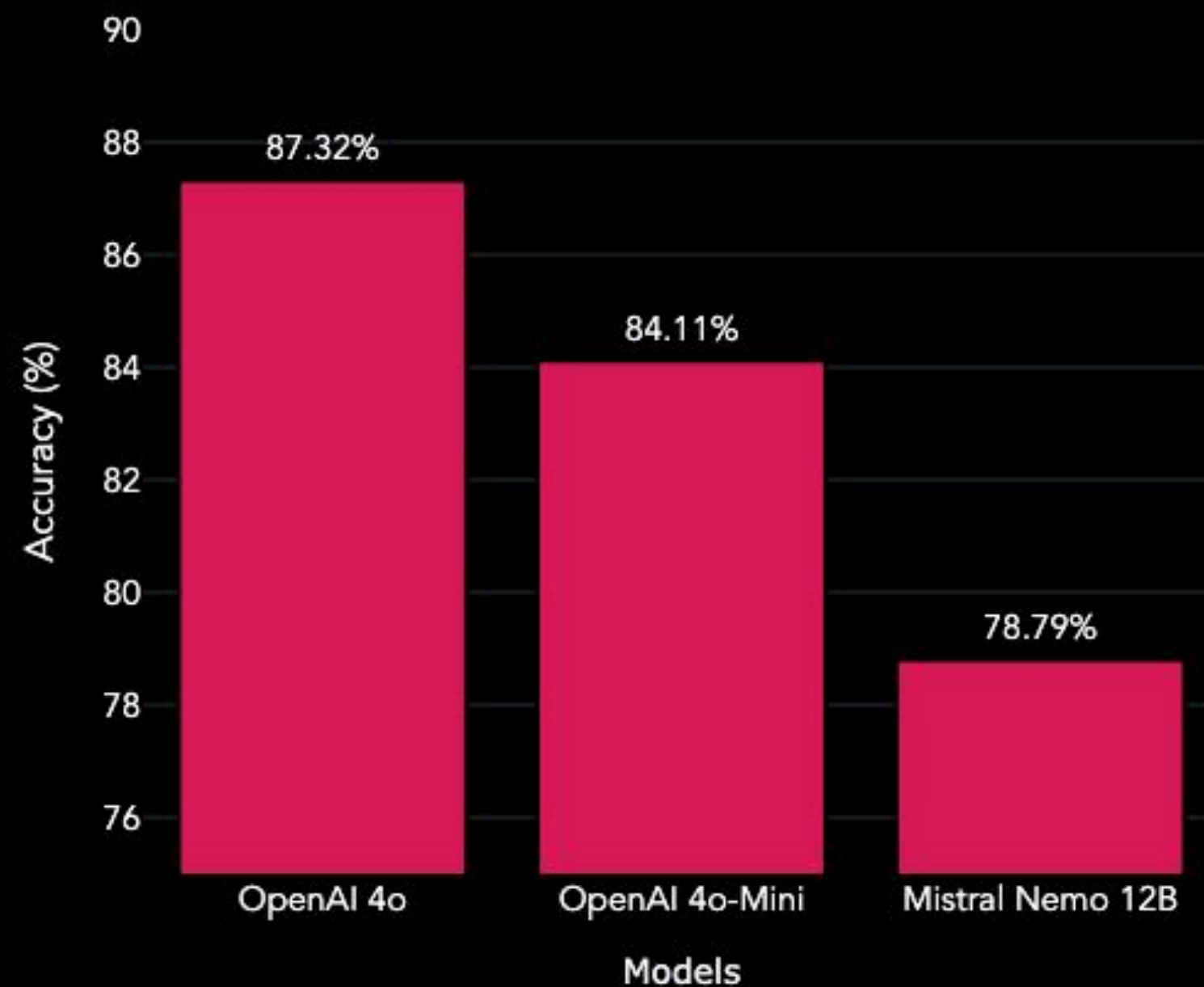


LLM Stage

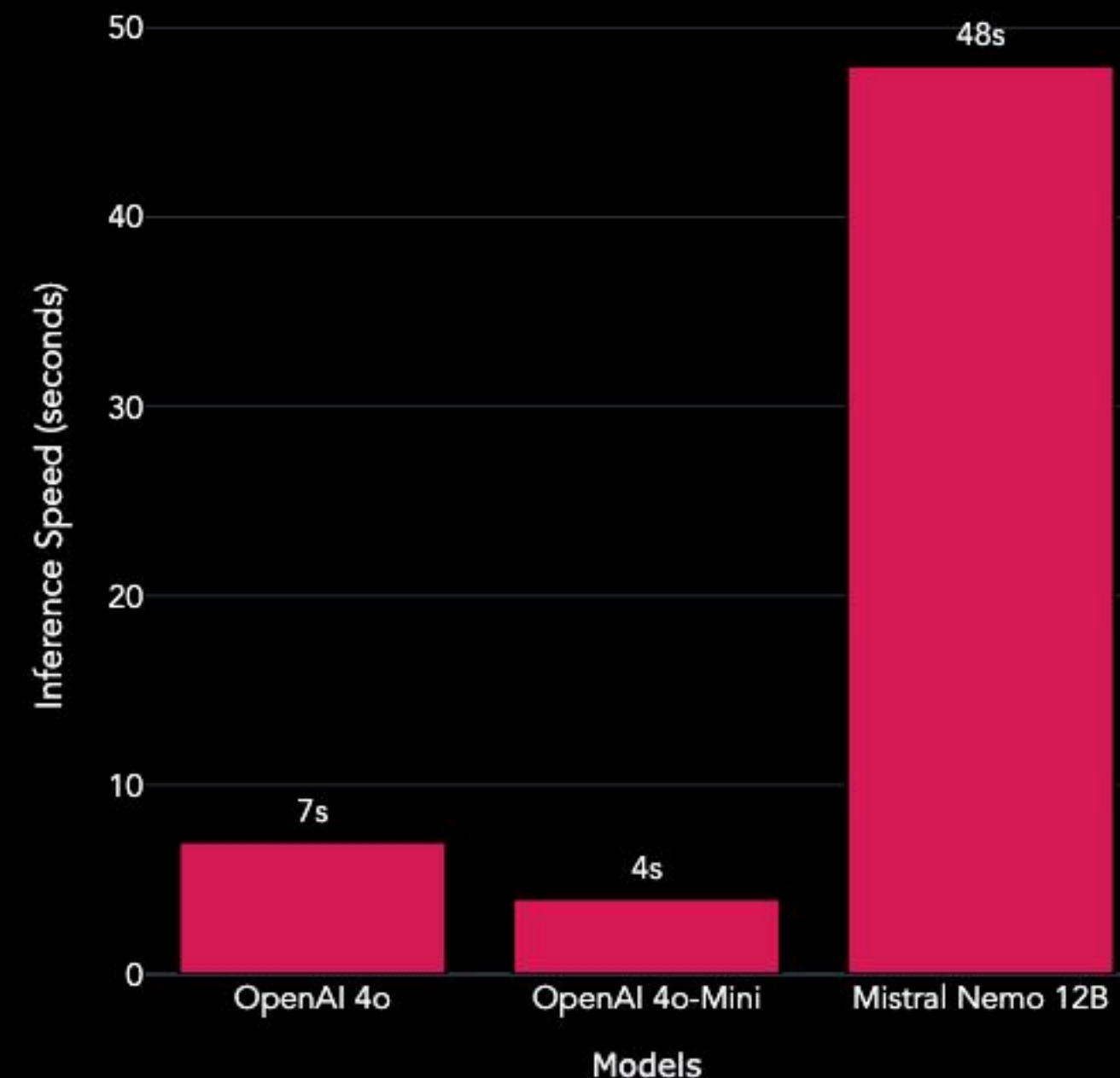
TESTING LARGE LANGUAGE MODELS

Open-Source vs Gated Models

ACCURACY



SPEED



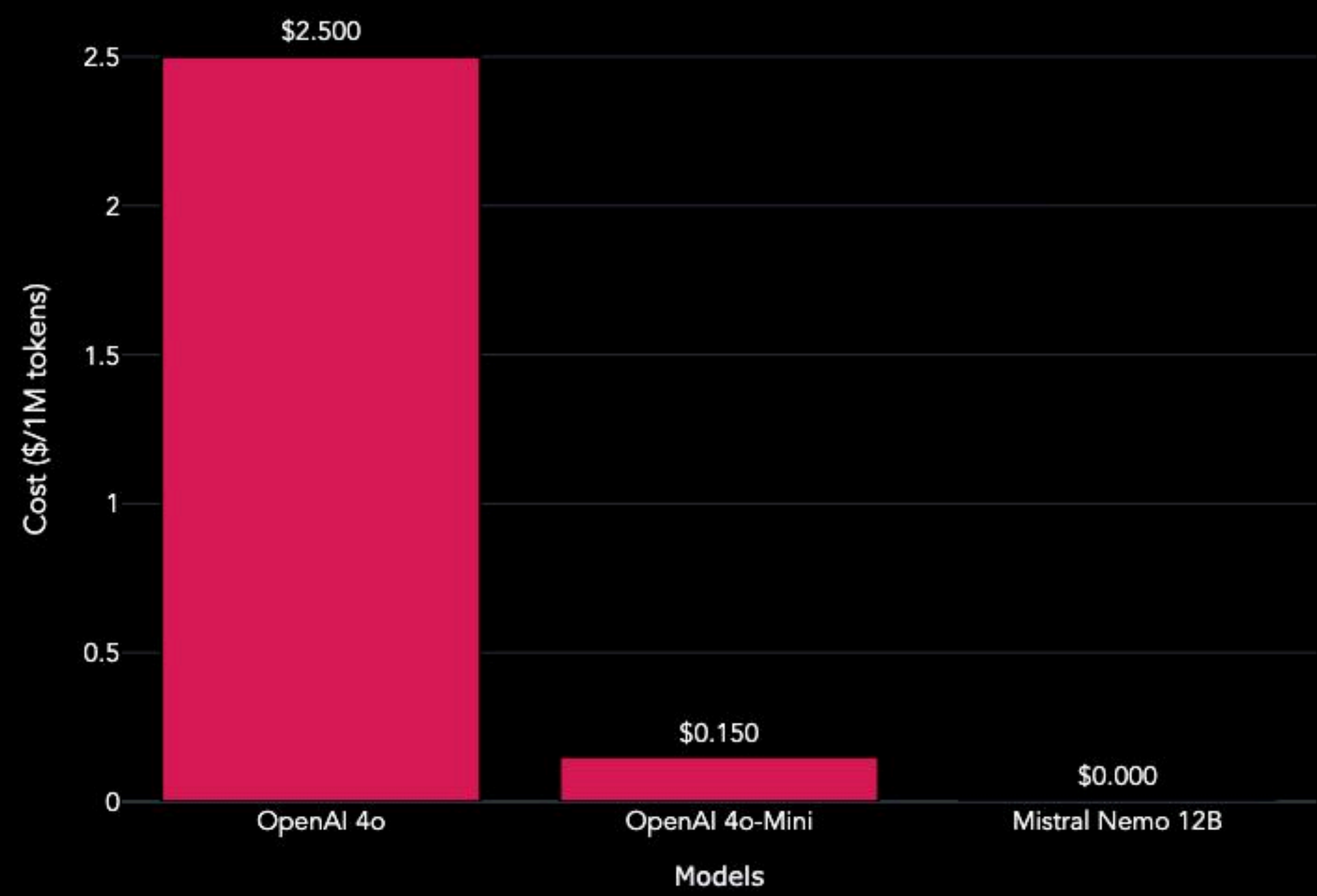
*2 V100 GPUs for Open-Source Models

Open-Source models are currently **not on-par** with OpenAI models for our task.

LLMS TESTED

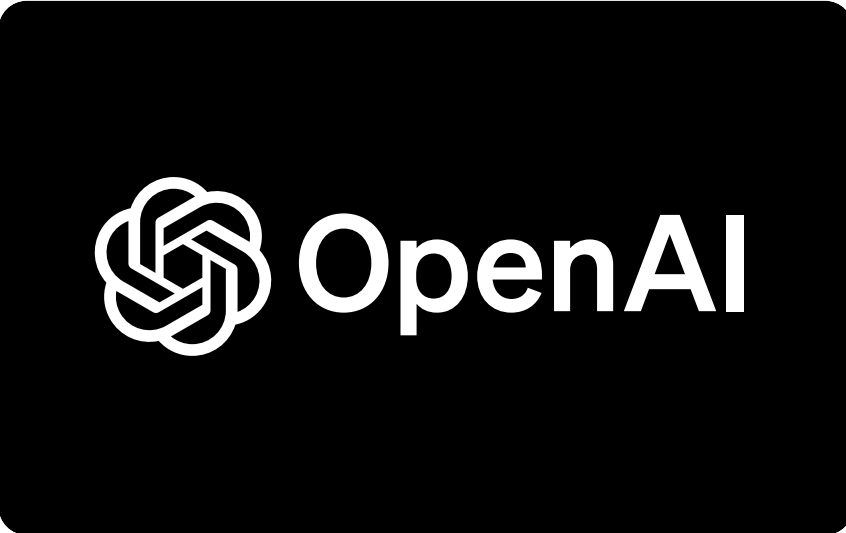


Open-Source vs Gated Models



As an initial start, we believe the cost is justified.
However, as we scale, we must consider otherwise.

LLMS TESTED



Model Validation

ACCURACY IN THE STRINGS WORLD

Levenshtein Distance

Counts the minimum number of single-character edits (insertions, deletions, or substitutions) needed to transform one string into the other.

We used **Tokenized** Levenshtein Distance.

1. Tokenize String - Ignores Weird Punctuation
2. Sort them Alphabetically - Ignores Order
3. Calculate Levenshtein Distance
4. Normalize by String Length - Get % Accuracy

Accuracy("Let's cook Grandma", "Let's cook, Grandma") = 100%

Insights taken from Model Validation were used to refine prompts.

IMPACT OF PROMPT ENGINEERING

72% Initial
with Original Prompt

87% Final
Prompt Changes from Validation Analysis

+15% Gain
More Potential.....

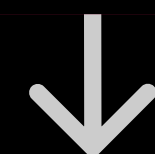
Final Pipeline

2 Page Input .TIF
xN



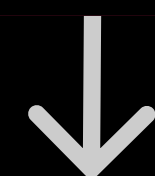
docTR OCR
Stage

(Resnet + VGG RCNN)

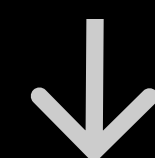


LLM Stage

OpenAI 4o



Sanity Checks



.CSV with Entity Extractions
Nx12

JSON
Output

JSON
Output

By The Numbers

Final Approach

87%

Accuracy

Levenshtein Distance

26_s

Per Document

OCR + API Call

.0006\$

Per Document

Average Doc is 2500 tokens

1773\$

For Entire Corpus

280,000 Documents

Next Steps

The Future

Scale

Scale to thousands of documents. Will understand more model nuances and refine prompts accordingly.

Go Open-Source

Llama 3.3 just pulled up. Put that to the test.

Publish

Goal has been to make the dataset and code public. So, once internal testing is done, publish.

Thanks.

Thanks.

Thanks.

Thanks.

Thanks.

Thanks.