# Artificial Intelligence (AI) & PDF Document Processing

Henry Sal
President Computing System Innovations
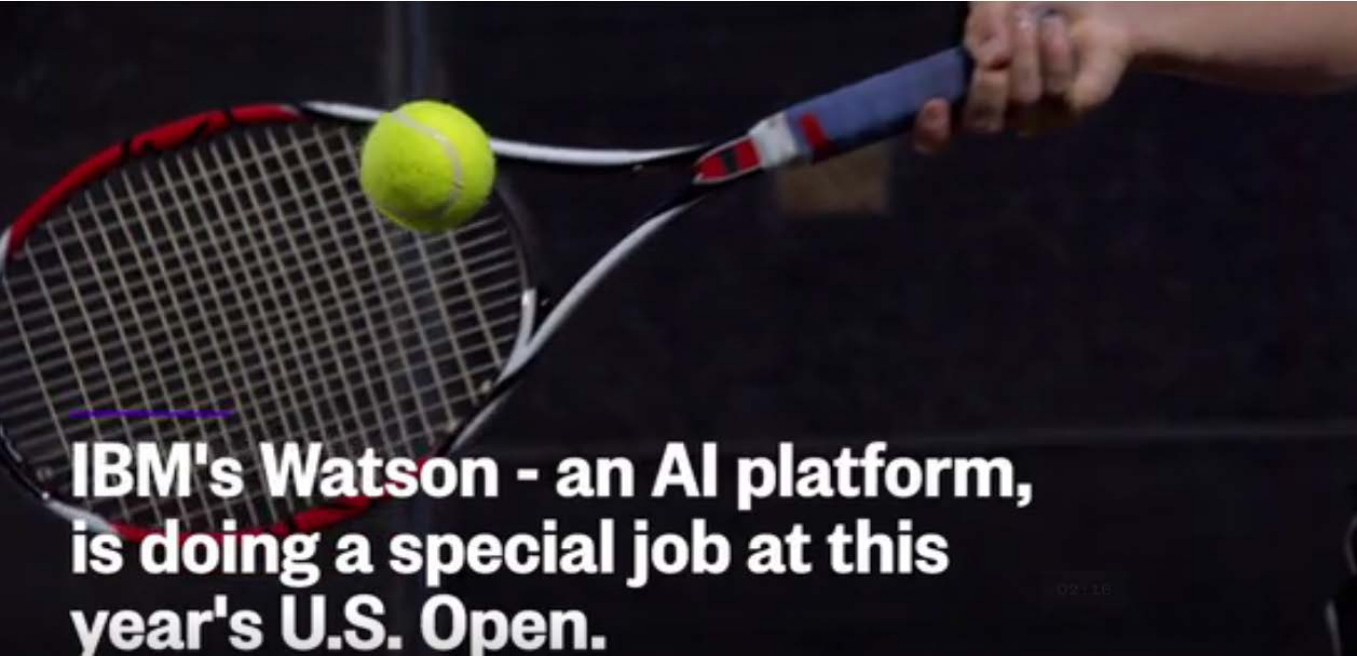
# Todays Agenda

- What is artificial intelligence (AI)?

- How does machine learning work?

- What is good machine learning?

- Practical AI PDF applications

- Accuracies & Exceptions

# Watson serves up Cognitive Highlights at the US Open

Cognitive

IBM's Watson - an AI platform, is doing a special job at this year's U.S. Open.

## COGNITIVE HIGHLIGHT OF THE DAY FACTORS

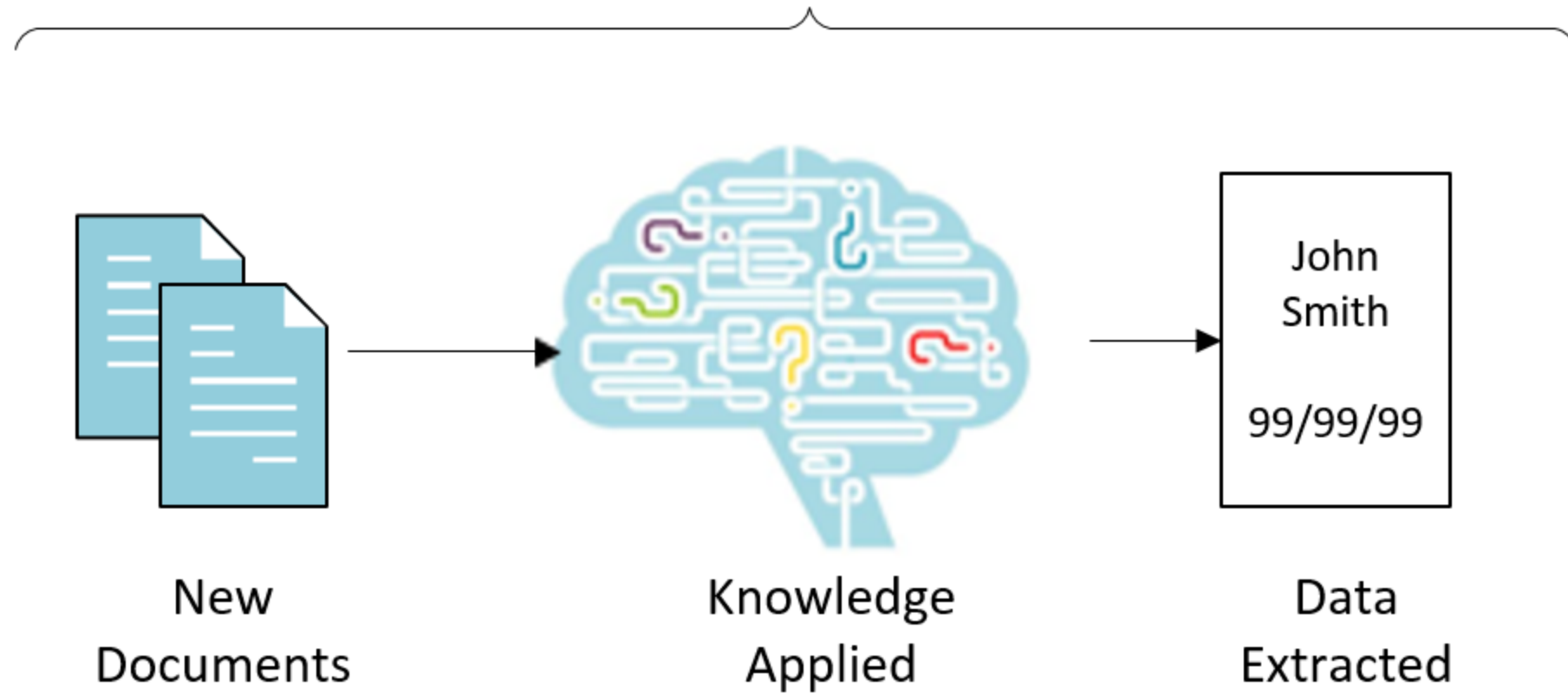| MATCH DATA | CROWD CHEERING | PLAYER GESTURE | OVERALL EXCITEMENT |
|---|---|---|---|
| .40 | 1 | .75 | .83 |

US OPEN | With Watson

CSI

# Types of AI

- **General AI** – simulated human intelligence, fabulous machines that have all our senses and more.  Deep learning, AlphaGo, Watson, etc…

- **Narrow AI** – Technologies that are able to perform specific tasks as well as, or better than, we humans can
  - Electronic stock trading
  - Facial recognition
  - Self-driving vehicles
  - Unstructured document analysis

- **Machine learning** – state-of-the art approach to narrow AI. Given example data, machines teach themselves and then make predictions within specific domains.

  - **Supervised learning** – provide initial examples of data that machine algorithms train on and learn how to perform that task.

  - **Online learning** – automatic improvement of learned knowledge via real world usage.
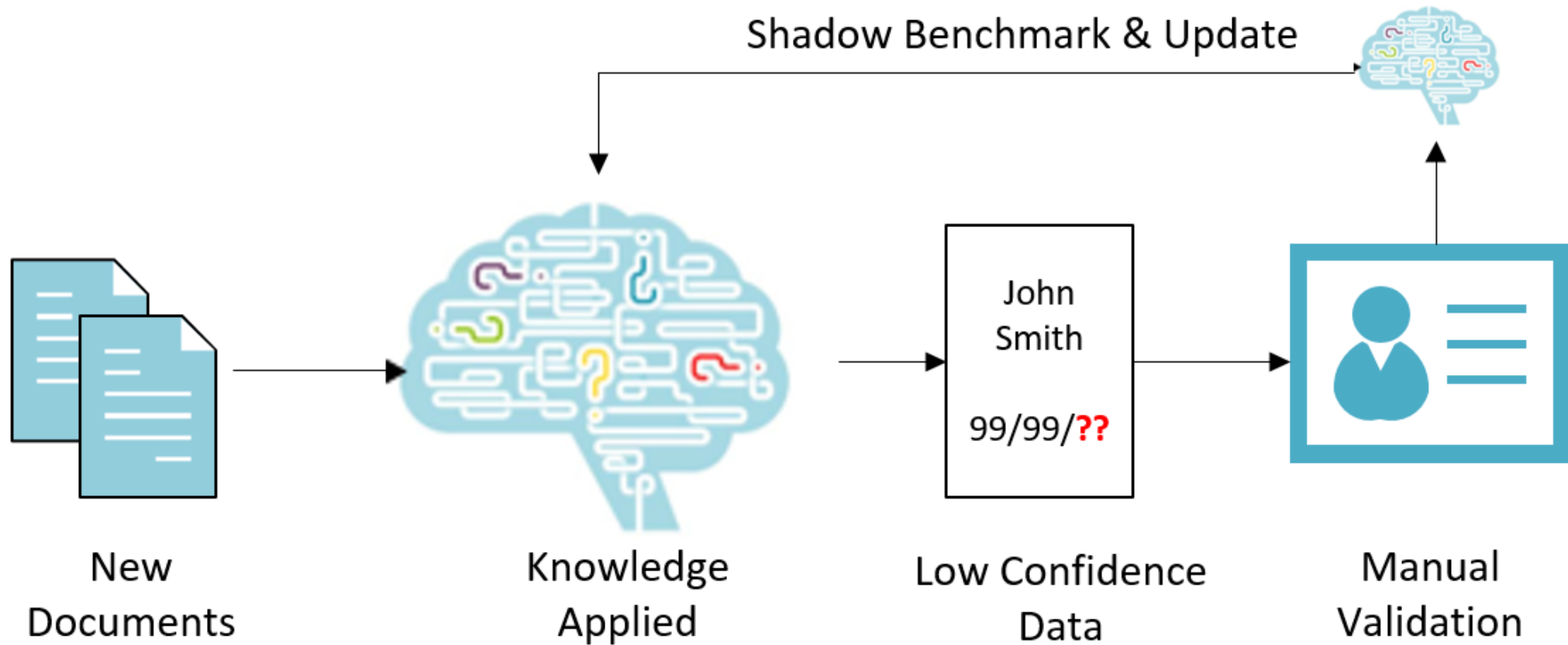
# Finding Text

# The Old Way
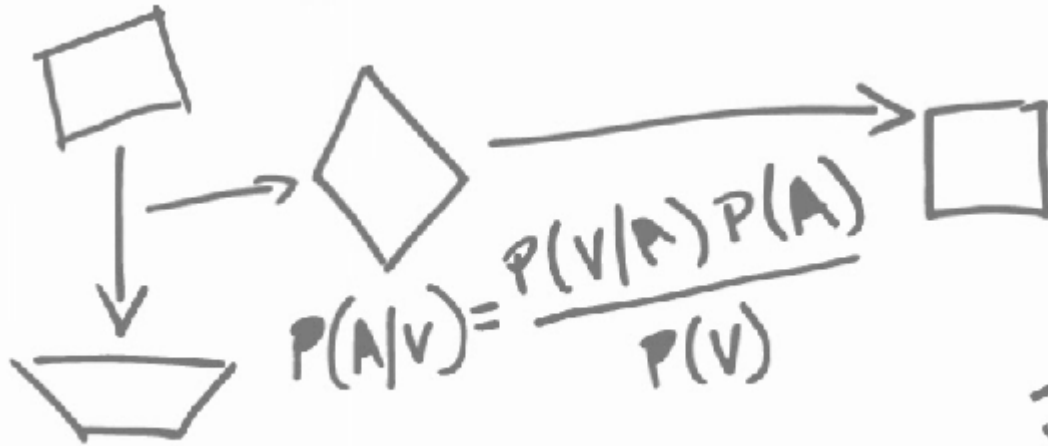
**Regular Expressions Handprint SSN**

([C][A]|[F][L]|[V][A]|<K california 90>|<K florida 95>|<K virginia 90>)[\ ]?[0-9]{3}[\- ][0-9]{2}[=-][0-9]{4}|

([#:\.\ ]([0-9]{3}[\ ]?[=\*\-][\ ]?[0-9ABbCDGHIiJKkLOQRSUZXx&]{2}[\ ]?[=\*\-][\ ]?[0-9ABbCDGHIiJKkLOQRSUZXx&]{4})[,;\.\ ])|
(^([0-9]{3}[\ ]?[=\*\-][\ ]?[0-9ABbCDGHIiJKkLOQRSUZXx&]{2}[\ ]?[=\*\-][\ ]?[0-9ABbCDGHIiJKkLOQRSUZXx&]{4})[,;\.\ ])|
([#:\.\ ]([0-9]{3}[\ ]?[=\*\-][\ ]?[0-9ABbCDGHIiJKkLOQRSUZXx&]{2}[\ ]?[=\*\-][\ ]?[0-9ABbCDGHIiJKkLOQRSUZXx&]{4})^)|
(^([0-9]{3}[\ ]?[=\*\-][\ ]?[0-9ABbCDGHIiJKkLOQRSUZXx&]{2}[\ ]?[=\*\-][\ ]?[0-9ABbCDGHIiJKkLOQRSUZXx&]{4})^)|
([#:\.\ ]([0-9][0-9ABbCDGHIiJKkLOQRSUZXx&]{2}[\ ]?[=\*\-][\ ]?[0-9ABbCDGHIiJKkLOQRSUZXx&]{2}[\ ]?[=\*\-][\ ]?[0-9]{4})[,;\.\ ])|
(^([0-9][0-9ABbCDGHIiJKkLOQRSUZXx&]{2}[\ ]?[=\*\-][\ ]?[0-9ABbCDGHIiJKkLOQRSUZXx&]{2}[\ ]?[=\*\-][\ ]?[0-9]{4})[,;\.\ ])|
([#:\.\ ]([0-9][0-9ABbCDGHIiJKkLOQRSUZXx&]{2}[\ ]?[=\*\-][\ ]?[0-9ABbCDGHIiJKkLOQRSUZXx&]{2}[\ ]?[=\*\-][\ ]?[0-9]{4})^)|
(^([0-9][0-9ABbCDGHIiJKkLOQRSUZXx&]{2}[\ ]?[=\*\-][\ ]?[0-9ABbCDGHIiJKkLOQRSUZXx&]{2}[\ ]?[=\*\-][\ ]?[0-9]{4})^)|
([#:\.\ ]([0-9][0-9ABbCDGHIiJKkLOQRSUZXx&]{2}[=\-\ ][0-9ABbCDGHIiJKkLOQRSUZXx&]{2}[=\-\ ][0-9]{4})[,;\.\ ])|
(^([0-9][0-9ABbCDGHIiJKkLOQRSUZXx&]{2}[=\-\ ][0-9ABbCDGHIiJKkLOQRSUZXx&]{2}[=\-\ ][0-9]{4})[,;\.\ ])|
([#:\.\ ]([0-9][0-9ABbCDGHIiJKkLOQRSUZXx&]{2}[=\-\ ][0-9ABbCDGHIiJKkLOQRSUZXx&]{2}[=\-\ ][0-9]{4})^)|
(^([0-9][0-9ABbCDGHIiJKkLOQRSUZXx&]{2}[=\-\ ][0-9ABbCDGHIiJKkLOQRSUZXx&]{2}[=\-\ ][0-9]{4})^)|
([#:\.\ ]([0-9][0-9ABbCDGHIiJKkLOQRSUZXx&]{2}[=\-\.][0-9ABbCDGHIiJKkLOQRSUZXx&]{2}[=\-\.][0-9]{4})[,;\.\ ])|
(^([0-9][0-9ABbCDGHIiJKkLOQRSUZXx&]{2}[=\-\.][0-9ABbCDGHIiJKkLOQRSUZXx&]{2}[=\-\.][0-9]{4})[,;\.\ ])|
([#:\.\ ]([0-9][0-9ABbCDGHIiJKkLOQRSUZXx&]{2}[=\-\.][0-9ABbCDGHIiJKkLOQRSUZXx&]{2}[=\-\.][0-9]{4})^)|
(^([0-9][0-9ABbCDGHIiJKkLOQRSUZXx&]{2}[=\-\.][0-9ABbCDGHIiJKkLOQRSUZXx&]{2}[=\-\.][0-9]{4})^)|
([Xx]{2,3}[=\-\.][Xx]{2}[=\-\.][0-9]{4})|
([#:\.\ ]([0-9]{3}[=\-\ ][0-9]{2}[=\-\ ][KkXx*]{4})[,;\.\ ])|
(^([0-9]{3}[=\-\ ][0-9]{2}[=\-\ ][KkXx*]{4})[,;\.\ ])|
([#:\.\ ]([0-9]{3}[=\-\ ][0-9]{2}[=\-\ ][KkXx*]{4})^)|
(^([0-9]{3}[=\-\ ][0-9]{2}[=\-\ ][KkXx*]{4})^)|
<K claim 95>[\ ]<K number 95>[:]?[\ ]{1,2}([0-9]{3}[\-]?[0-9]{2}[\-][0-9]{4}|[0-9]{3}[\-][0-9]{2}[\-]?[0-9]{4})

Online Learning

Shadow Benchmark & Update

John
Smith

99/99/**??**

New
Documents

Knowledge
Applied

Low Confidence
Data

Manual
Validation

CSI

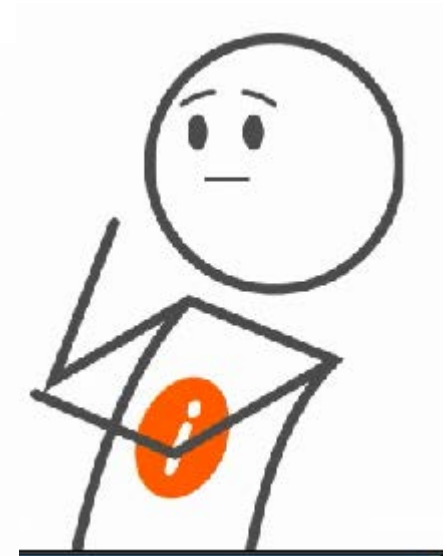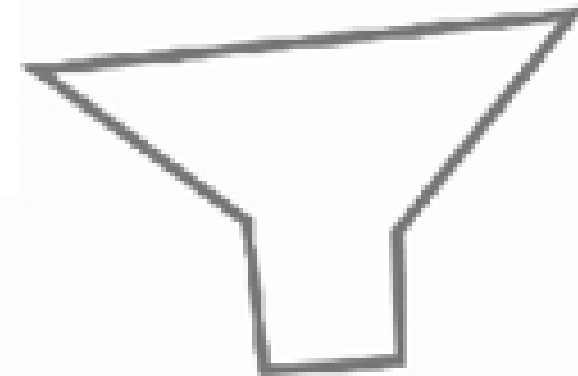# Content Locators, Role Determination, Conditional Random Fields (CRFs), Natural Language Processing (NLP), Parts of Speech Tagging, etc...



$$P(A|V) = \frac{P(V|A)\,P(A)}{P(V)}$$

$$J(\bullet) = \frac{1}{2M} \sum_{i=1}^{m} (h_\theta(x^{(i)})$$

$$\bullet = \frac{1}{1 + \exp\left[-\left(\beta_0 + \sum_{r} \beta_r\right)\right]}$$

# Visible features

- Format
- Relative location
- What's around it (that matters)

# Invisible features

- Semantics (parts of speech)

Filing<VVG> #<UNC> 13325244<UNC> Electronically<AVo> Filed<VVD> 05/06/2014<CRD> 01<CRD>:<PUN>46<CRD>:<PUN>28<CRD> PM<AVo>

IN<PRP> THE<ATo> CIRCUIT<NN1> COURT<NN1> OF<PRF> THE<ATo> NINTH<ORD> JUDICIAL<AJo> CIRCUIT<NN1>,<PUN> IN<PRP> AND<CJC> FOR<PRP> ORANGE<AJo> COUNTY<NN1>,<PUN> FLORIDA<NPo>

CASE<NN1> NO<ITJ>.<PUN>:<PUN> 2014-CA-001234-0<PUQ>

DBK<VDB> MANAGEMENT<NN1>,<PUN> LLC<VHG>,<PUN> a<ATo> Florida<NPo> limited<AJo> liability<NN1> company<NN1>,<PUN>

Plaintiff<NPo>,<PUN>

v<CRD>.<PUN>

HRKALOVIC<AJo> HOMES<NN2>,<PUN> LLC<ITJ>,<PUN> a<ATo> Florida<NPo> limited<AJo> liability<NN1> company<NN1>,<PUN>

Defendant<NN1>.<PUN>

NOTICE<NN1> OF<PRF> HEARING<VVG> PLEASE<AVo> TAKE<VVB> NOTICE<NN1> that<CJT> on<PRP> the<ATo> 15<CRD>th<NNo> day<NN1> of<PRF> July<NPo>,<PUN> 2014<CRD>,<PUN> at<PRP> 2<NNo>:<PUN>00<CRD> p<ZZo>.<PUN>m<ZZo>,<PUN>,<PUN> or<CJC> as<CJS> soon<AVo> thereafter<AVo> as<CJS> counsel<NN1> can<VMo> be<VBI> heard<VVN>,<PUN> the<ATo> undersigned<AJo> counsel<NN1> for<PRP> Defendant<NN1>,<PUN> HRKALOVIC<AJo> HOMES<NN2>,<PUN> LLC<ITJ>,<PUN> a<ATo> Florida<NPo> limited<AJo> liability<NN1> company<NN1>,<PUN> will<VMo> call<VVI> up<AVP> for<PRP> hearing<VVG> before<PRP> the<ATo> Honorable<AJo> Donald<NPo> E<ZZo>.<PUN> Grincewicz<NPo>,<PUN> in<PRP> Courtroom<NN1> 18B<CRD>,<PUN> of<PRF> the<ATo> Orange<NN1> County<NN1> Courthouse<NN1>,<PUN> 425<NNo> North<NN1> Orange<AJo> Avenue<NN1>,<PUN> Florida<NPo> 32801<UNC>,<PUN> the<ATo> following<AJo> matter<NN1>:<PUN>

MOTION<NN1> TO<TOo> DISMISS<VVI> COUNTS<PNX> II<CRD> AND<CJC> IV<CRD> OF<PRF> PLAINTIFFS<NN2> COMPLAINT<NN1> OR<CJC> FOR<PRP> A<ATo> MORE<AVo> DEFINITE<AJo> STATEMENT<NN1> AS<CJS> TO<PRP> COUNT<VVI> II<CRD>

PLEASE<AVo> GOVERN<VVB> YOURSELVES<PNX> ACCORDINGLY<AVo>.<PUN>

i<PNP>4t<ORD>-<PNI> Respectfully<AVo> submitted<VVN> this<DTo> _<NN1>£<NNo>_<VHG>day<NN1>

# Example Knowledge



Extraction Pattern Configuration Dialog - EventDate

Trained Patterns

| Phrase | Role ▼ | Context | Backward | Conflict |
|---|---|---|---|---|
| Motion | Event Date | Nort-West from entity | ☑ | 0 % |
| NOTICE OF HEARING | Event Date | North from entity | ☑ | 0 % |
| has been set as follows | Event Date | North from entity | ☑ | 0 % |
| vs | Event Date | North, Nort-West from entity | ☑ | 0 % |
| AMENDED NOTICE OF HEARING | Event Date | North from entity | ☑ | 0 % |
| CERTIFICATE OF SERVICE | Event Date | South from entity | ☑ | 9 % |
| Plaintiff | Event Date | North from entity | ☑ | 6 % |
| Time | Event Date | Suffix phrase, directly right of entity | ☑ | 0 % |
| hearing on | Event Date | Prefix phrase, directly left of entity | ☑ | 0 % |
| on the | Event Date | Prefix phrase, directly left of entity | ☑ | 0 % |
| PLEASE BE GOVERNED ACCORDI | Event Date | South from entity | ☑ | 0 % |
| Judge | Event Date | South-West from entity | ☑ | 0 % |

Items displayed: 44

Close

# What is good machine learning?

- Highly accurate results

- Doesn't require enormous volume of examples to train on

- Learns from its mistakes, not a static fragile system

- Flexible! Users can easily establish new data items without software development

- Not a black box. Analysts can review, debug, and refine knowledge

- Most steps performed by software (i.e. initial tagging, model refinement, online learning, etc…)

# Practical Application to PDF Documents

- Automatically locate and extract (or redact) data

- Auto tag PDF documents for accessibility

- Auto separate and bookmark embedded PDFs

- Transform unstructured content to structured output(s)

- Eliminate human document review & data entry

- Create 24x7x365 "lights out" document workflows

# Accuracies & Exceptions

- OCR engines are pattern recognition engines
  (Machine learning from the '70s)

- OCR full page reads 90%+
  (Image quality, skew, background, form dropouts, color, etc...)

- Field level accuracy 85% - 100%

- Transaction accuracy 85%+  — end user gold standard

- Digitally born PDFs accuracy game changer — no OCR needed

# Questions?

# Thank You!

[hsal@csisoft.com](mailto:hsal@csisoft.com)



PDF ■ DAY
WASHINGTON DC