## Collaborative PDF

Leonard Rosenthol, Sr. Principal Architect, PDF & Content Authenticity



#### Collaboration

#### Dictionary

Definitions from Oxford Languages · Learn more



#### noun

1. the action of working with someone to produce something.

"he wrote a book in collaboration with his son"

2. traitorous cooperation with an enemy.

"he faces charges of collaboration"

## Who (or what) are we looking to collaborate with?



# What sorts of Collaboration are we looking for?

- Viewing
- Comprehension
- Commenting & Annotating
- Editing
- Manipulation/Transformation



# What sorts of Collaboration are we looking for?

- Viewing
- Comprehension
- Commenting & Annotating
- Editing
- Manipulation/Transformation



## Comprehension



#### **Structural Semantics**

- The current PDF tagging functionality
  - aka the "Structure Tree"
  - How a document's content is constructed

- DParts Document Parts
  - Relationships of Pages w/associated "attributes" (aka DPM, DPart Metadata)



### Derivation (to HTML & beyond)

- PDF Association's Derivation Algorithm
  - https://pdfa.org/resource/deriving-html-from-pdf/
  - Tagged PDF -> HTML

- LLMs love HTML (& understand general semantics such as headings, lists & tables)
- HTML can be easily converted (with loss) to Markdown, for those cases where the LLM prefer that.



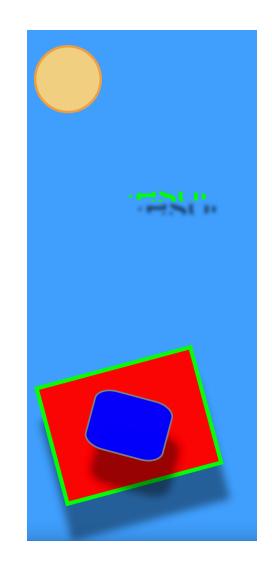
## What about the PDF-centric "goodness"?

- XMP Metadata (not just PDF, but still)
- Annotations (incl. hyperlinks)
- Form Fields (incl. Signatures)
- Associated/Embedded Files



#### Missing: Content Semantics

- Explicitly identify objects
  - Regular polygons with rounded corners
  - Text on a curve
- Explicitly identify object components
  - Shadow cast on an object
  - Effects (e.g., glow, blur)
- Explicitly identify grouping objects
  - Axis components on a graph inside a Figure structural element

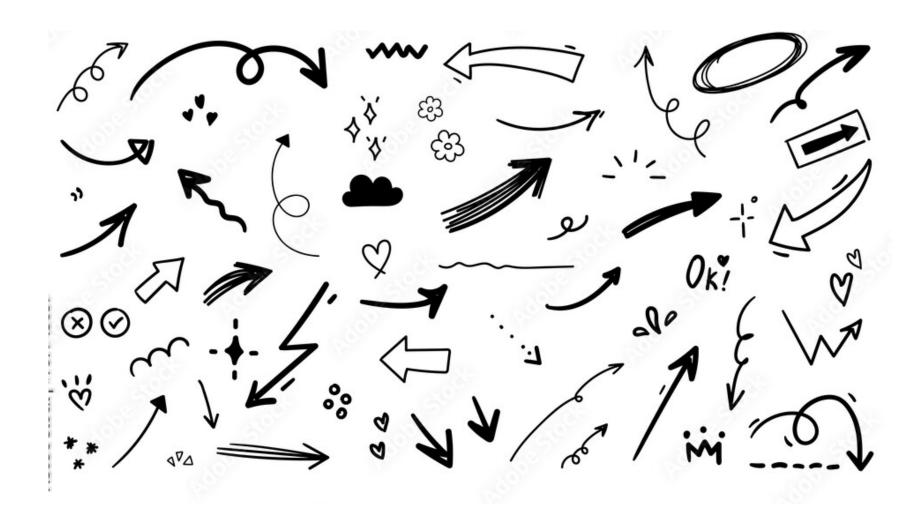


### How to express all this to AI?

- Find a "markup language" it may have already trained on
  - SVG works for many things (e.g., object grouping, effects, etc.), but not others (e.g., reflowable text)
  - But LOTS of stuff there isn't a single option out there
- Try a "data driven approach"
  - Create a JSON grammar that can describe content + attributes
  - Convince an existing LLM to use it
- Train your own model

Throw up your hands and run screaming

## Commenting & Annotating



#### **Annotations**

- Text Notes
- "Mark Up"
  - Underline, strikeout, highlight
  - Circle, Rect, Pen/Scribble
  - Callout, Dimensioning
- Hypertext Links
  - Inter-document, intra-document
  - URI's (web links)
- Stamps
- Rich Media (incl. 3D)



#### Information Age Insecurity

This is a simple text annotation

The Information Age is irrevocably altering the means by which the Government must approach the challenge of protecting its information. Protection no longer equates to placing documents in filing cabinets with strong combination locks. Instead, information vital to the security and continued prosperity of the United States resides in a series of increasingly interconnected classified and unclassified systems. The Commission believes that the findings and recommendations noted below provide policymakers the means to begin protecting information properly now and into the next century.

This is an era of extraordinary change not only in information technology, but also in the very waying which individues sommunion with one another. The Commission's goal is not the predict the floure that these rechlological changes will help mold. Rather, it is to better understand the nature of the new throats so that the Governmen with the full support of the private sector, car mitigate or between them.

At present, there exists what appears to be a growing gap between technological change and the human capacity to adapt to that change. The risk is that the Government will make bad decisions not because it has too little information, but rather because it has too much information about the wrong things. In such a rapid-paced and changing environment, it is only natural to fall back on old biases, protocols, and shortcuts. Convictions, as Nietzsche once noted, can be "more dangerous enemies of truth than lies."

#### Federal Government Information Security and the National Information Infrastructure

The information revolution, characterized by the growing convergence of computer and communications technologies, requires a fundamental rethinking of traditional approaches to safeguarding national security information. Those responsible for the

of national security face new, increasingly difficult challenges presented by ation of computer networks linked by telephone lines, cable, direct broad, and wireless communications, and by the replacement of the traditional nainframe by personal computers. In this new electronic world—the formation Infrastructure (NII)—best symbolized by the steadily growing met, it is not clear what responsibility the Federal Government has to protect ucture that stores, carries, and transmits nearly all of the Government's I and classified information.

thin the United States is only one portion of the Global Information Infra-GII) that connects public and private computer networks around the world. It is described by the Global Information of the WII, which is a leadership position in protecting the NII,







### Markup Annotation "Extras" (1.5+)

- IRT "In Response To"
  - Supports hierarchical relationships incl. grouping for annotations
  - Used for threaded replies

#### State

- Reviewed (i.e., approved, completed, rejected, etc.)
- Marked



1

#### **Bhavin Kapoor**

May 31

@JayGhoshKhemka - Shared Document links being opened in Acrobat Web - 40.6M



Jay Ghosh Khemka May 31

Hi @HarshpreetKaur: Can you please confirm if 40.6M is app launch or doc opens? Can you please share the logic or AA workspace for this?

Jay Ghosh Khemka

May 31

@AnubhaJain: FYI

Reply

Page 2



#### **Bhavin Kapoor**

May 31

Is this also last one month count Or is this 2 week count?

II.



May 3

@VarinderSaini - Can you please help to get this checked.

The first 3 items as MAU, is this last item also tracking MAU or is this 2 weeks data

Reply



## What's Missing?

- Document level annotations
- Page level annotations
- Reactions/"Liking"/Acknowledging
  - Emojis (which may requires newer Unicode versions!)
- Multi vs. single threading



#### **AcroForms**

- A PDF file contains a single AcroForm, though that form may contain any number of fields located on any number of pages.
- There a number of predefined field types
  - Button (checkbox/radio/push)
  - Text
  - Choice (popup, combo or list)
  - Digital Signatures
- Fields can be typed (integer, string, boolean) and marked read-only
- Fields can be "calculated"



#### How to express all this to AI?

- No common grammars/languages can express everything that our annotations support.
  - What info is most important and does that change based on use case?

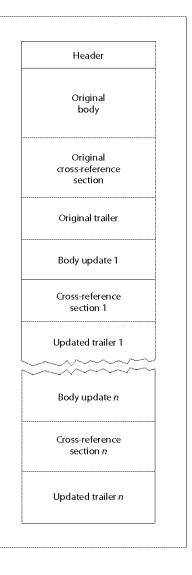
- Same for Forms especially when you want to include business logic, formatting, etc.
  - The Forms.Next work of the PDF Association a modern JSON grammar could help here!



# Editing

### Incremental Update

- Modifications are written to the end of the file, leaving the original data intact
- A new xref table is written containing the new/modified data, and a link back to the old xref.
- Since original data is still present, support for multiple undos across save boundaries could be supported.



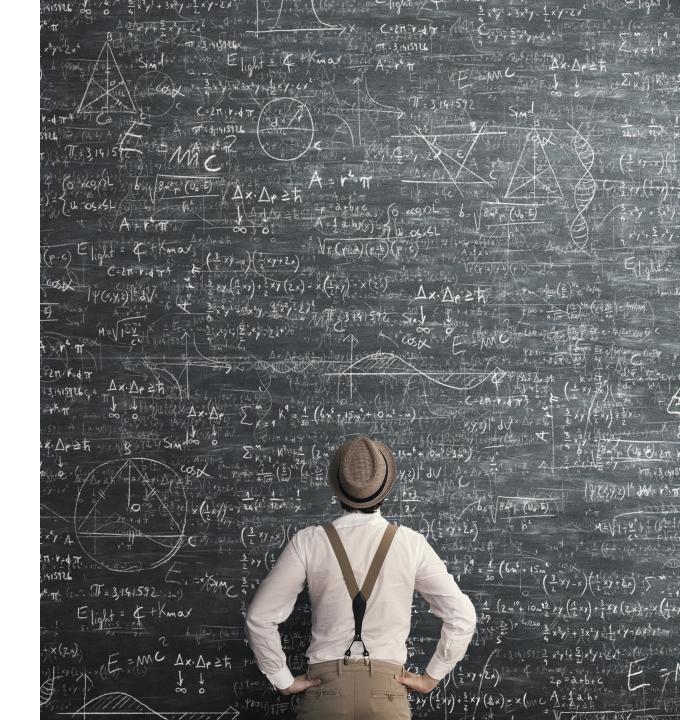
### Incremental update assumes single user

- Heavy weight for simple changes
  - One character change on a page impacts many objects (at least a full content stream)
- Not (PII) safe across removal operations
  - Delete Page, Redact, etc.
- Single threaded no forks or branches



#### Annotations + Editing = ?!?!?

 Having things like links and markup separated from the content significantly complicates allowing it to be (a) extracted and (b) edited.





## Content Provenance



#### ISO 22144: Content Credentials

An open standard for storing and accessing <u>cryptographically</u>
 <u>verifiable and tamper-evident information</u> whose
 trustworthiness can be assessed based on a <u>defined trust model</u>.



#### A look inside of Content Credentials



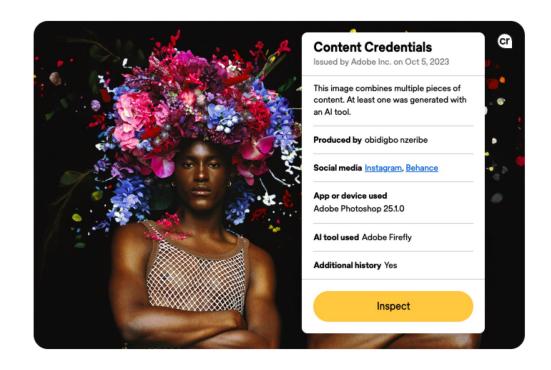


#### Functionality around Generative Al

Identifying assets & regions of interest that have been created/modified by Al

Information about the "recipe"

Allow creators to label content "Do Not Train"



#### ISO 32008: PDF Extensions to Content Credentials

- Extending 22144 & the C2PA spec for more PDF-isms
- New Actions
  - Over 20 proposals to discuss this week!
- Better support for Certifying Signatures
- Handling C2PA Manifests in encrypted PDFs
- Handling of C2PA Manifest with Protected Wrappers (AU).
- Guidance on image placement & transcoding

# ISO 21617 (JPEG Trust): Extensions to Content Credentials for identity, authorship, ownership & digital rights

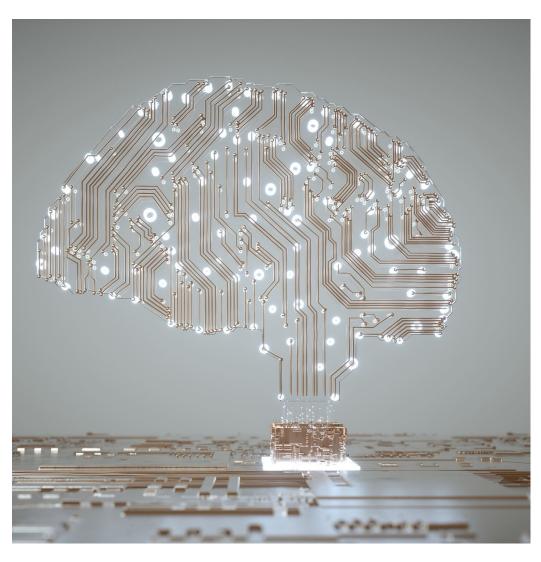
- Assert human and organizational identities using standard models
  - X.509
  - Verifiable Credentials
  - Social Media accounts
  - Well known ID systems
- Describe role(s) involved in the production of the asset
- Provide rights declarations (ODRL)
  about the asset & its content





# Agents

#### Al Agents Transforming PDF Comprehension and Processing



- Al agents analyze PDF structure for enhanced content understanding.
- They extract text, images, and annotations efficiently.
- Enable automated commenting, editing, and metadata management.
- Improve accessibility through semantic tagging and format conversion.
- Facilitate content provenance and collaboration in PDFs.

# Conclusion

#### Conclusion

- PDFs are the "next frontier" for content consumption
- But WE UNDERSTAND they are so much more than content
- How do we help the world get there?

- It's not just about improving the file format (we need that too) –
  but about finding ways to helps users better use it!
  - Both creation and consumption

Myou J04" HUOU you gout jank thank nk nk you thank thank you that nank thank thank thank

#### Questions



# Adobe

- In a world where users expect to collaborate when authoring or editing content, where AI agents are out there performing actions on behalf of their users (or even themselves!), what does it mean for PDF a file format designed to allow a single user to do their work?
- I'll talk about what the modern collaborative era means to PDF from the file format, to attribution and provenance. That is, assuming that I don't send my Agent to present on my behalf!