# Deriving HTML from PDF

How reusable content in pdf could be

Roman Toda, Normex

electronic document CONFERENCE

# History, Evolution of PDF

- Electronic representation of paper. Based on graphical model, accurate representation (screen and print)
  - Adopted word **rendering**
- Marked content
  - Introducing semantic into the content
- Tagged PDF
  - First time used different presentation model for assistive technology

# Problems 2019

- PDF is essential part of the web – but is it really?

- HTML users hate pdf
  - How to use pdfs on mobiles
  - In browsers

- HTML developers hate pdf
  - Can't control the user experience
  - Can't access/navigate pdf content

Roman Toda  NORMEX

# Derivation

- Is it possible to deterministically interpret PDF content in the web based environment?

- Yes. **Tagged PDF** is the best way for capturing author's intent
  - provisions in ISO 32000-2 are rich enough to be unambiguously interpreted in html = **Derived**

- **Derivation** our new word for **rendering**

Roman Toda NORMEX

# Join us

- Download and read the publication

- PDF Association

- Next-Generation PDF working group

- Test our MVP
  - https://github.com/Normex/PDF-Derivation
  - Share your files
  - Discuss techniques
  - Give feedback

Roman Toda  NORMEX

# What is Tagged PDF

- Structure tree
  - Marked content
  - Forms fields, links, annotations
- Attributes
- Classes
- Associated files
- Actions

Roman Toda  NORMEX

# What is Tagged PDF

# What is Tagged PDF

# Standard structure element types

- Chapter 14.8.4 in ISO 32000-2

    - Annot, Artifact, Aside, Caption, Document, DocumentFragment, Div, Em, FEnote, Figure, Form, Formula, H1..Hn, L, Lbl, LBody, LI, Link, P, RB, RP, RT, Ruby, Span, Strong, Sub, Table, TBody, TD, TFoot, TH  THead, Title, TR, Warichu, WT, WP

- Some sound like html tags: Div, Span, P, H1, LI

- Some are harder to imagine: Annot

- RoleMap may be involved

# Text

# Simple structure - Text

# Table
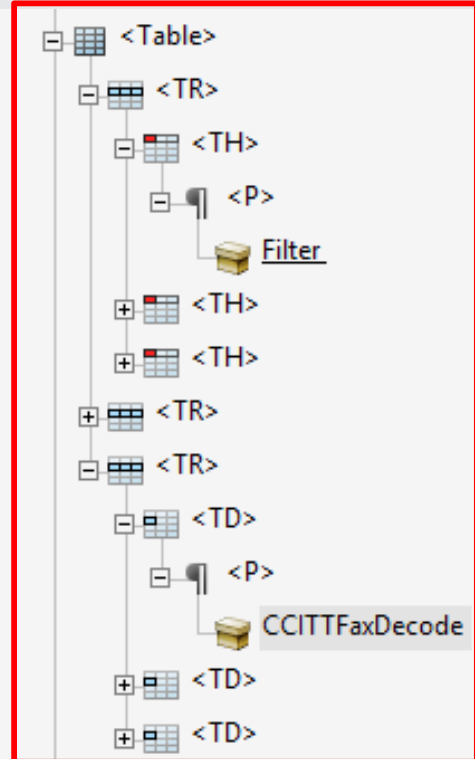
# Table - derivation



```
<table>
    <tr>
        <th>
            <p lang="EN-US">
                <span>Filter </span>
            </p>
        </th>
        <!--.....-->
    </tr>
    <!--.....-->
    <tr>
        <td>
            <p lang="EN-US">
                <span>CCITTFaxDecode </span>
            </p>
        </td>
        <!--.....-->
    </tr>
</table>
```
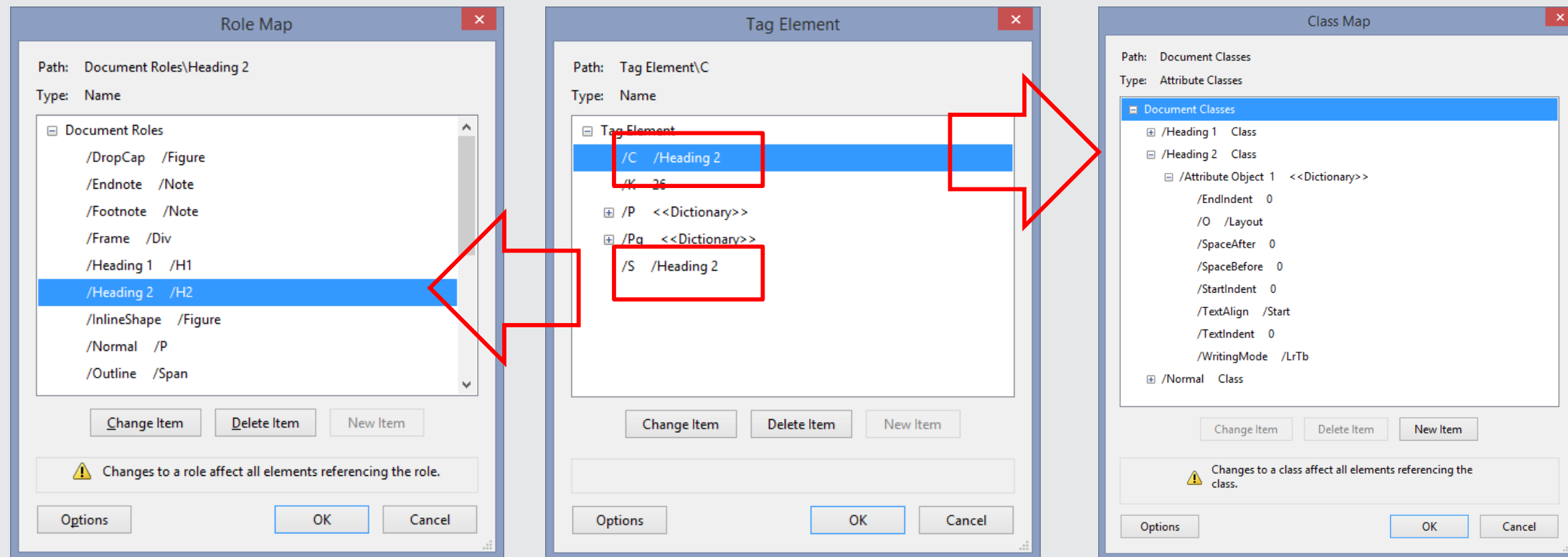
# Table - derivation



```
<table>
    <tr>
        <th>
            <p lang="EN-US">
                <span>Filter </span>
            </p>
        </th>
        <!--.....-->
    </tr>
    <!--.....-->
    <tr>
        <t
```

```
        </
        <!
    </tr>
</table>
```

**Raster images**

Raster images in PDF (called Image XObjects) are represented by dictionaries with an associated stream.

| Filter | Extension | Description |
|---|---|---|
| DCTDecode | .jpg | a lossy filter based on the JPEG standard |
| CCITTFaxDecode | .tiff | a lossless bi-level (black/white) filter based on the Group 3 or Group 4 CCITT (ITU-T) fax compression standard defined in ITU-T T.4 and T.6 |

# Attributes – general concept

- Each tag could have attributes
- Different practice is using ClassMap
- Role mapping takes place

# Attribute owners

**Table 376: Standard structure attribute owners**

| Owner value for the attribute object's O entry | Description |
| --- | --- |
| Layout | Attributes governing the layout of content |
| List | Attributes governing the numbering of lists |
| PrintField | Attributes governing Form structure elements for non-interactive form fields |
| Table | Attributes governing the organization of cells in tables |
| Artifact | Attributes governing Artifact structure elements |
| XML-1.00 | Additional attributes governing translation to XML, version 1.00 |
| HTML-3.20 | Additional attributes governing translation to HTML, version 3.20 |
| HTML-4.01 | Additional attributes governing translation to HTML, version 4.0 |
| HTML-5.00 | Additional attributes governing translation to HTML, version 5.0 |
| OEB-1.00 | Additional attributes governing translation to OEB (Open eBook), version 1.0 |
| RTF-1.05 | Additional attributes governing translation to Microsoft Rich Text Format, version 1.05 |
| CSS-1.00 | Additional attributes governing translation to a format using CSS, version 1.00 |
| CSS-2.00 | Additional attributes governing translation to a format using CSS, version 2.00 |
| CSS-3.00 | Additional attributes governing translation to a format using CSS, version 3.00 |
| RDFa-1.10 | Additional attributes governing translation to a format using RDFa version 1.1 |

# Attribute owners

**Table 376: Standard structure attribute owners**

| Owner value for the attribute object's O entry | Description |
|---|---|
| Layout | Attributes governing the layout of content |
| List | Attributes governing the numbering of lists |
| PrintField | Attributes governing Form structure elements for non-interactive form fields |
| Table | Attributes governing the organization of cells in tables |
| Artifact | Attributes governing Artifact structure elements |
| XML-1.00 | Additional attributes governing translation to XML, version 1.00 |
| HTML-3.20 | Additional attributes governing translation to HTML, version 3.20 |
| HTML-4.01 | Additional attributes governing translation to HTML, version 4.0 |
| HTML-5.00 | Additional attributes governing translation to HTML, version 5.0 |
| OEB-1.00 | Additional attributes governing translation to OEB (Open eBook), version 1.0 |
| RTF-1.05 | Additional attributes governing translation to Microsoft Rich Text Format, version 1.05 |
| CSS-1.00 | Additional attributes governing translation to a format using CSS, version 1.00 |
| CSS-2.00 | Additional attributes governing translation to a format using CSS, version 2.00 |
| CSS-3.00 | Additional attributes governing translation to a format using CSS, version 3.00 |
| RDFa-1.10 | Additional attributes governing translation to a format using RDFa version 1.1 |

Roman Toda    NORMEX

# Attribute usage



- Layout
- Table
- List
- Can use the power of CSS

# Attribute usage

The Portable

## Imaging model

The basic design of how gra[...]
except for the use of transpa[...]

The current transformation [...]

- The clipping path
- The color space
- The alpha constant,

## Raster images

Raster images in PDF (called[...]
associated stream.

| Filter | Ext |
|---|---|
| DCTDecode | .jpg |
| CCITTFaxDecode | .tif |

**Attributes**

Path:    Attribute Objects\

Type:    Dictionary

☐ Attribute Objects

☐ /Attribute Object 1    <<Dictionary>>

/O    /CSS-3.00

/font-size    (36px)

/text-align    (center)

/color    (blue)

/background-color    (yellow)

/word-spacing    (10px)

[ Change Item ] [ Delete Item ] [ New Item ]

```
<h1 lang="EN-US"
style="background-color: yellow;
color: blue;
font-size: 36px;
text-align: center;
word-spacing: 10px; ">
<span>Imaging model </span>
</h1>
```

# Attribute usage



```
<h1 lang="EN-US"
style="background-color: yellow;
color: blue;
font-size: 36px;
text-align: center;
word-spacing: 10px; ">
<span>Imaging model </span>
</h1>
```

Roman Toda   NORMEX

# Associated files – css, js, svg, mathml

# Associated files – css, js, svg, mathml

# Namespaces – sample with html



Roman Toda  NORMEX

# Namespaces – sample with html



Roman Toda   NORMEX

# Namespaces – sample with html

# Forms sample

# Forms sample

# Valid html



Roman Toda  NORMEX

# Summary

- World of single representation is over. Today we need to read out loud our pdfs, but what's going to be tomorrow?

- It's time to stop producing bad pdfs

- It's time to give authors full control over pdf interpretation

- Author decides instead of tool decides

- Instead of saying there is only one way to interpret PDF we will be saying that there is **always deterministic** way of interpreting pdf

Roman Toda  NORMEX

# Resources

- [https://github.com/Normex/PDF-Derivation](https://github.com/Normex/PDF-Derivation)

- Sample files

  - Styling

  - Associated files

  - Forms

  - Interactive

  - Fail cases

- Implementation (commandline, GUI)

# Resources



Roman Toda · NORMEX

**?**

Thanks !

Roman Toda, Normex

https://github.com/Normex/PDF-Derivation