# Deriving HTML from PDF lessons learned

Implementation challenges in reusing PDF content

PDF Days Online 2021

- 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
  - PDF/UA
  - First time used different presentation model for assistive technology

- Problems – PDF in the web, in mobiles
  - Can't control UX
  - Interactivity, navigation
  - Responsivness
- Solution – Deriving HTML from PDF
  - Algorithm that produces conforming HTML from a tagged PDF
  - Derivation = HTML rendering
  - Author controls the output not the processor

# HTML Derivation - concept

- "converting" pdf structure elements to html tags

- Structure tree

  - Marked content

  - links, annotations, form fiels

- Attributes (Layout, CSS, HTML) via ClassMap too

- Associated files

- Actions

- Rolemap, Namespaces

# Valid HTML

| child \| / parent --> | | StructTree Root (div) | Document (div) | Document Fragment (div) | Part (div) | Div (div) | Aside (aside) | Title (p) | Sub (span) | P (p) | Hn (h1-h6/H7=p) | H (h1-h6) | Lbl (label/div) | Em (em) | Strong (strong) | Span (span) | Link (a) | Annot (??) | Form (form) | Ruby (ruby) | RB (rb) | RT (rt) | RP (rp) | Warichu (span) | WT (span) | WP (span) | FENote (p) | L (ul/ol) | LI (li) | LBody (div) | Table (table) | TR (tr) | TH (th) | TD (td) | THead (thead) | TBody (tbody) | TFoot (tfoot) | Caption (caption/figcaption) | Figure (figure) | Formula (figure) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Document | div | 0..1 | 0..n | 0..n | 0..n | 0..n | 0..n | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| DocumentFragment | div | | 0..n | 0..n | 0..n | 0..n | 0..n | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Part | div | | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | | | | | | | | | | | | | | | | | | | 0..n | | 0..n | | | | | | | | | | 0..n | 0..n | 0..n |
| Div | div | | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | | | | | 0..n | | | 0..n | 0..n | 0..n | | | | | | | | | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n |
| Aside | aside | | 0..n | 0..n | 0..n | 0..n | | 0..n | | | | | | | | | | | | | | | | | | | 0..n | | 0..n | | | | | | | | | 0..n | 0..n | 0..n |
| P | p | | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | | | | | | | | | | | | | | | | | | | 0..n | | 0..n | | | 0..n | 0..n | | | | | 0..n | 0..n | 0..n |
| Hn | h1-h6/H7=p | | 0..n | 0..n | 0..n | 0..n | 0..n | | | | | | | | | | | | | | | | | | | | | | 0..n | | | 0..n | 0..n | | | | 0..n | 0..n | 0..n |
| H | h1-h6 | | 0..n | 0..n | 0..n | 0..n | 0..n | | | | | | | | | | | | | | | | | | | | | | 0..n | | | 0..n | 0..n | | | | 0..n | 0..n | 0..n |
| Title | p | | 0..n | 0..n | 0..n | 0..n | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Sub | span | | | | 0..n | | | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | | | 0..n | 0..n | 0..n | | 0..n | 0..n | 0..n | | | 0..n | | | | | | | | 0..n | 0..n | 0..n |
| Lbl | label/div | | | | 0..n | | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | | | | | | | | 0..n | 0..n | | | | 0..n | 0..n | 0..n | | | | 0..n | 0..n | 0..n |
| Em | em | | | | | | | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | | | 0..n | 0..n | 0..n | | 0..n | 0..n | 0..n | | | 0..n | | | | | | | | 0..n | 0..n | 0..n |
| Strong | strong | | | | | | | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | | | 0..n | 0..n | 0..n | | 0..n | 0..n | 0..n | | | 0..n | | | | | | | | 0..n | 0..n | 0..n |
| Span | span | | | | | | | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | | | 0..n | 0..n | 0..n | | 0..n | 0..n | 0..n | | | 0..n | | | | | | | | 0..n | 0..n | 0..n |
| Link | a | | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | | | 0..n | 0..n | 0..n | | 0..n | 0..n | 0..n | | | 0..n | | | | | | | | 0..n | 0..n | 0..n |
| Annot | ?? | | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | | | 0..n | 0..n | 0..n | | 0..n | 0..n | 0..n | | | 0..n | | | | | | | | 0..n | 0..n | 0..n |
| Form | form | | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | | | | 0..n | 0..n | 0..n | | 0..n | 0..n | 0..n | | | 0..n | | | | | | | | 0..n | 0..n | 0..n |
| Ruby | ruby | | | | | | | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | | | | | | | | | 0..n | | | 0..n | | | | | | | | 0..n | 0..n | 0..n |
| RB | rb | | | | | | | | | | | | | | | | | | | [a] | | | | | | | | | | | | | | | | | | | | |
| RT | rt | | | | | | | | | | | | | | | | | | | [a] | | | | | | | | | | | | | | | | | | | | |
| RP | rp | | | | | | | | | | | | | | | | | | | [a] | | | | | | | | | | | | | | | | | | | | |
| Warichu | span | | | | | | | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | | | | | | | | | 0..n | | | 0..n | | | 0..n | 0..n | | | | 0..n | 0..n | 0..n |
| WT | span | | | | | | | | | | | | | | | | | | | | | | | [b] | | | | | | | | | | | | | | | | |
| WP | span | | | | | | | | | | | | | | | | | | | | | | | [b] | | | | | | | | | | | | | | | | |
| FENote | p | | | | 0..n | 0..n | 0..n | | | 0..n | | | | | | | | | | | | | | | | | | | | 0..n | 0..n | | 0..n | 0..n | | | | 0..n | 0..n | 0..n |
| L | ul/ol | | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | | | | | | | | | | | | | | | | | 0..n | 0..n | | 0..n | | | 0..n | 0..n | | | | 0..n | 0..n | 0..n |
| LI | li | | | | | | | | | | | | | | | | | | | | | | | | | | | | 0..n | | | | | | | | | | | |
| LBody | div | | | | | | | | | | | | | | | | | | | | | | | | | | | | | 0..n | | | | | | | | | | |
| Table | table | | | | | | | 0..n | | | | | | | | | | | | | | | | | | | 0..n | | | 0..n | | | 0..n | 0..n | | | 0..n | 0..n | 0..n |
| TR | tr | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | 0..n | | | 0..n | 0..n | 0..n | | | |
| TH | th | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | 0..n | | | | | | | | |
| TD | td | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | 0..n | | | | | | | | |
| THead | thead | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | 0..1 | | | | | | | | |
| TBody | tbody | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | 0..n | | | | | | | | |
| TFoot | tfoot | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | 0..1 | | | | | | | | |
| Caption | ption/figcaption | | | | 0..1 | 0..1 | 0..1 | 0..1 | | | | | | | | | | | | 0..1 | | | | | | | | | | | 0..1 | 0..1 | 0..1 | | | | | | 0..1 | 0..1 |
| Figure | figure | | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | | | | | | | | | 0..n | | | 0..n | | | 0..n | 0..n | | | | 0..n | 0..n | 0..n |
| Formula | figure | | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | 0..n | | | | | | | | | 0..n | | | 0..n | | | 0..n | 0..n | | | | 0..n | 0..n | 0..n |

# Resources

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

- Sample files
  - Styling
  - Associated files
  - Forms
  - Interactive
  - Fail cases

- Implementation (commandline, GUI)

# Problems - general

- Bad tagging (wrong attributes, content vs. struct tree)

- No PDF 2.0 support yet

- PDF/UA isn't Well Tagged = what is Well Tagged?

    - reusable pdfs also accessible ?

- No PDF/UA-2 yet

- Export doesn't produce Tagged PDFs (changing !!)

- No tools for manual tagging

# Problem - technical

- Elements (TOC)

- Attributes (Layout vs. CSS)

- Links across pages

- Forms, Javascript

- Annotations

# HTML vs. PDF

Deriving HTML from PDF

WT PDF → derived HTML

# HTML vs. PDF

Deriving HTML from PDF

WT PDF → derived HTML

PDF ← HTML

# HTML vs. PDF

# HTML vs. PDF

Deriving HTML from PDF

WT PDF → derived HTML

PDF

PDF/UA

Tagged PDF

accessibility

HTML

# HTML vs. PDF

Deriving HTML from PDF

WT PDF → derived HTML

?? ??

HTML

NORMEX

Foxit

# HTML 2 PDF

- Not a one to one conversion.
  - <div>, <span> → Div, Span
  - <b>, <i> → Span and loosing information about tag
  - <ul>, <ol> → L (ListNumbering) and we need to add Lbl that doesn't exist in html
  - <header>, <time> → loosing semantic, no attributes to represent such
- Complex problems
  - Interactive elements
  - Javascript
  - css

# HTML 2 PDF

- As many Layout attributes as possible

# HTML 2 PDF

- As many Layout attributes as possible

- 2.0 allows adding CSS, HTML attributes

# HTML 2 PDF

- As many Layou[t]
  possible

- 2.0 allows add[...]
  attributes

# HTML 2 PDF

- As many Layout attributes as possible

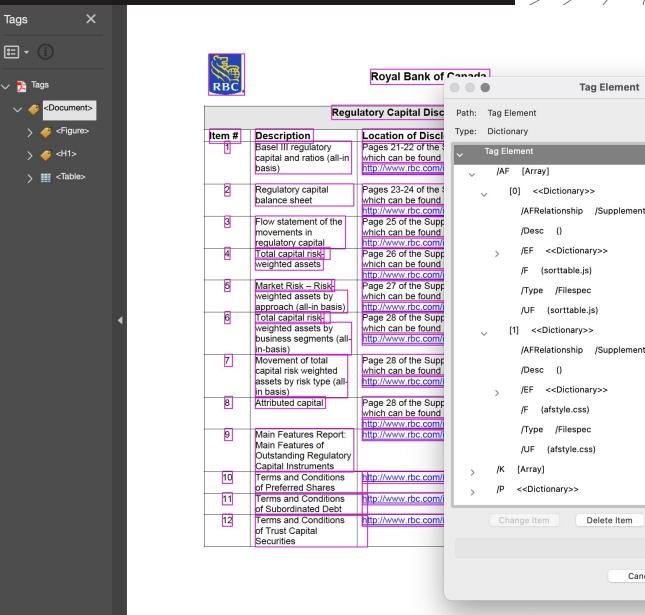- 2.0 allows adding CSS, HTML attributes

- URL binding

# HTML 2 PDF

- As many Layout attributes as possible

- 2.0 allows adding CSS, HTML attributes

- URL binding

- Associated files

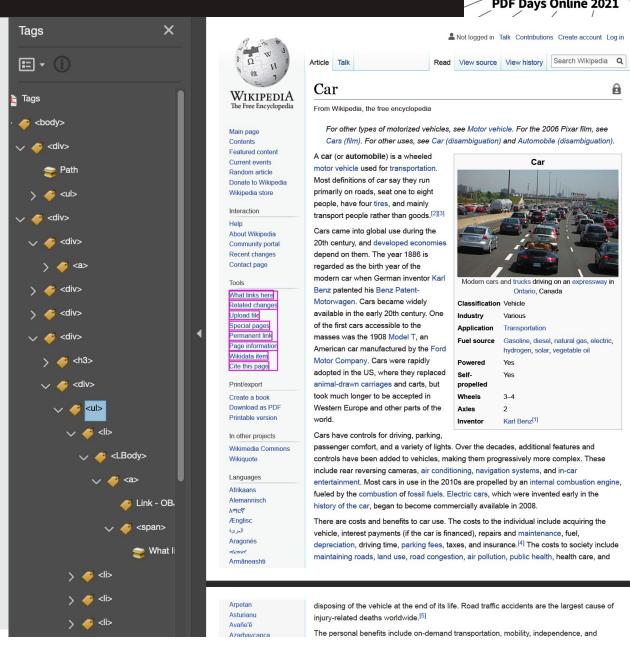# HTML 2 PDF

- As many Layout attributes as possible

- 2.0 allows adding CSS, HTML attributes

- URL binding

- Associated files

- HTML Namespace

# HTML 2 PDF

- As many Layout attributes as possible

- 2.0 allows adding CSS, HTML attributes

- URL binding

- Associated files

- HTML Namespace

- Forms

# Lessons learned - problems

- HTML – small, PDF - big

- Invalid html

- Incompatible tagsets

- complex html, css, js frameworks, media queries etc..

- Interactivity (forms, fieldsets - import/export)

- Online vs. offline (in html world we reference everything by URI)

- Hard to include multiple targets (UA, WT PDF – or mobile vs. desktop)

# Landscape change

- WT PDF technical working group

- Forms technical working group

- AnnexL update

- 1.7 and 2.0 Namespaces


- PDF/UA-2

# Join us

- Implementation
  - https://github.com/Normex/PDF-Derivation
- Samples
  - https://normex.github.io/PDF-Derivation/
- PDF Association TWG
  - Derivation TWG
  - WT PDF TWG
  - PDF/UA TWG

**?**

# Thanks !

Roman Toda, Normex

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

PDF association