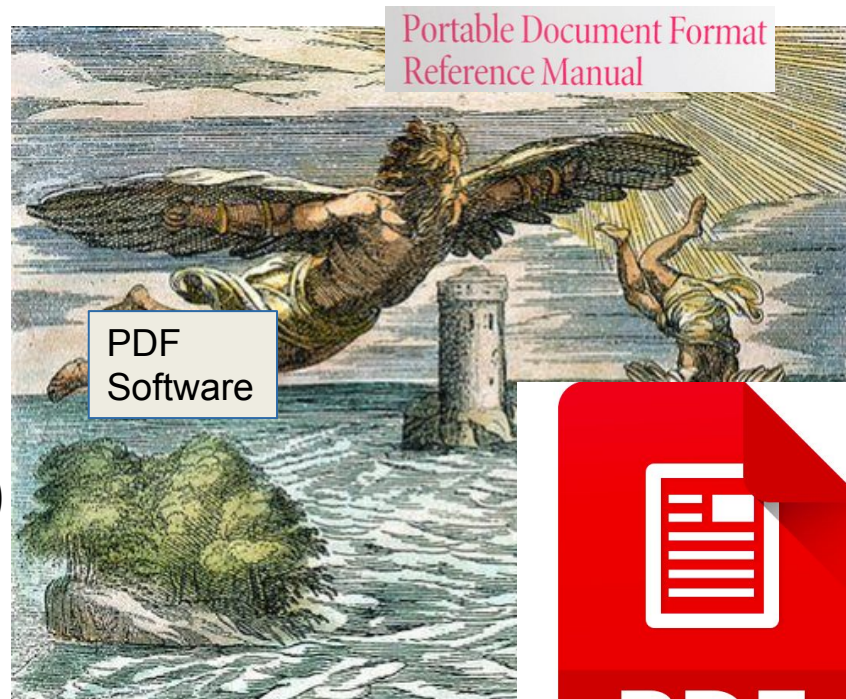**Research Report:**

# ICARUS: Understanding De Facto Formats By Way of Feathers and Wax

Sam Cowger, Yerim Lee, Nichole Schimanski, Mark Tullsen, **Walt Woods**, Richard Jones, EW Davis, William Harris, Trent Brunson, Carson Harmon, Bradford Larsen, Evan Sultanik
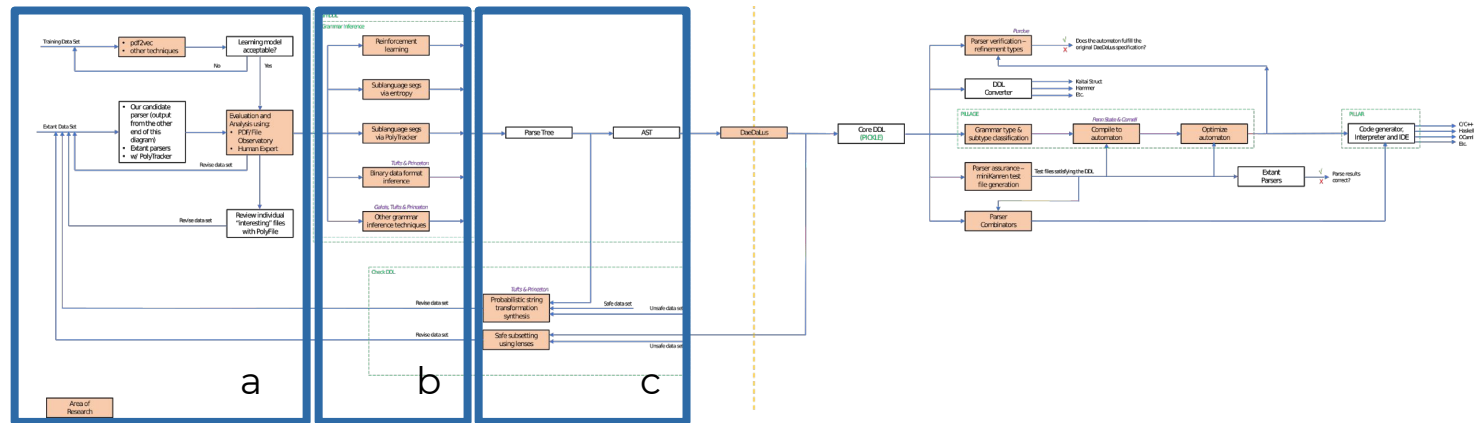
# De Facto Corpuses

- Widely adopted formats can expand beyond original specifications.

- Can affect data:
  - At rest (format ambiguity)
  - In use (parser vulnerabilities)
  - In transit (exfiltration)

Portable Document Format Reference Manual

PDF Software

PDF

# ICARUS Project

- Pipeline and tools focused on developing methodologies for:
  a. Discovering and describing de facto data formats;
  b. Identifying grammars within a de facto format; and
  c. Translating from a de facto format to a safe subset.

Galois SafeDocs Structural Overview

# ICARUS Project

- Pipeline and tools focused on developing methodologies for:
  a. Discovering and describing de facto data formats;
  b. Identifying grammars within a de facto format; and
  c. Translating from a de facto format to a safe subset.

Galois SafeDocs Structural Overview

# File Observatory

- Need to understand which collected files actually belong to de facto, vs malformed.

- Different parsers accept different files.

- Can learn a lot from stdout and stderr

Image from ImpulseCreative.com

# File Observatory

```
53        ^parsers_qpdf_Parser_WARNING: loop detected following xref tables
54        # This one is unreliable, and refers to previous lines in output ^pa
55        ^parsers_mutool_Parser_error: Unable to read ICC workflow
56        # This one often applies to xref issues; ^parsers_mutool_Parser_warn
57        ^parsers_caradoc_Parser_PDF error : Lexing error : integer error : i
58        ^parsers_mutool_Parser_error: malformed page tree
59
60 ˅    RejectedBad:
61        ^parsers_mutool_Parser_error: no objects found
62        ^parsers_caradoc_Parser_PDF error : Lexing error : unexpected charac
63        ^parsers_pdfinfo_ParserStruct_Syntax Error \(\): Illegal character <
64        ^parsers_pdftocairo_Parser_Syntax Error \(\): Illegal character '}'
65        ^parsers_pdftocairo_Parser_Syntax Error
66        ^parsers_qpdf_Parser_WARNING: operation for dictionary attempted on
```

**^parsers_pdfinfo_ParserStruct_Syntax Error: I**

**SafeWarnings:**
  **^parsers_qpdf_Parser_WARNING: unknown token w**
  **^parsers_pdfid_Parser_/URI**

```
77 ˅    UnsafeWarnings:
78        ^parsers_qpdf_Parser_WARNING: loop detected following xref tables
79        ^parsers_pdfid_Parser_/OpenAction
80        ^parsers_pdfinfo_ParserNorm_Syntax Error \(\): Dictionary key must be
81        ^parsers_mutool_Parser_warning: non-page object in page tree
82
83 ˅  outputs:
84        # Standard output status -- If a PDF passes filter S1, it will be "val
85        # otherwise "rejected".
86 ˅    status:
87        "valid" is !(RejectedBad | RejectedAmbiguousBad | ValidWarningsXrefR
88        "rejected" else
89
```
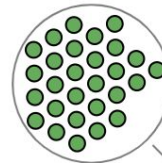
Special filters

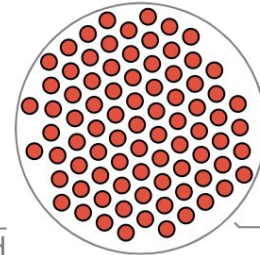parsers_caradoc_Parser alphanumeric      parsers_pdfid_Parser alphanu

parsers_pdfinfo_ParserStruct struct      parsers_pdfinfo_ParserMeta not

Line-item filters allowed (click to reject)

✓ parsers_pdftocairo_Parser_Syntax Error (): Dictionary key must be a name obje
✓ parsers_pdftocairo_Parser_Syntax Error (): Illegal character '>'

valid                                    rejected

REPROCESS DECISIONS      DOWNLOAD DECISIONS      REPROCESS DB ERRORS

# ICARUS Project

- Pipeline and tools focused on developing methodologies for:
  a. Discovering and describing de facto data formats;
  b. Identifying grammars within a de facto format; and
  c. Translating from a de facto format to a safe subset.

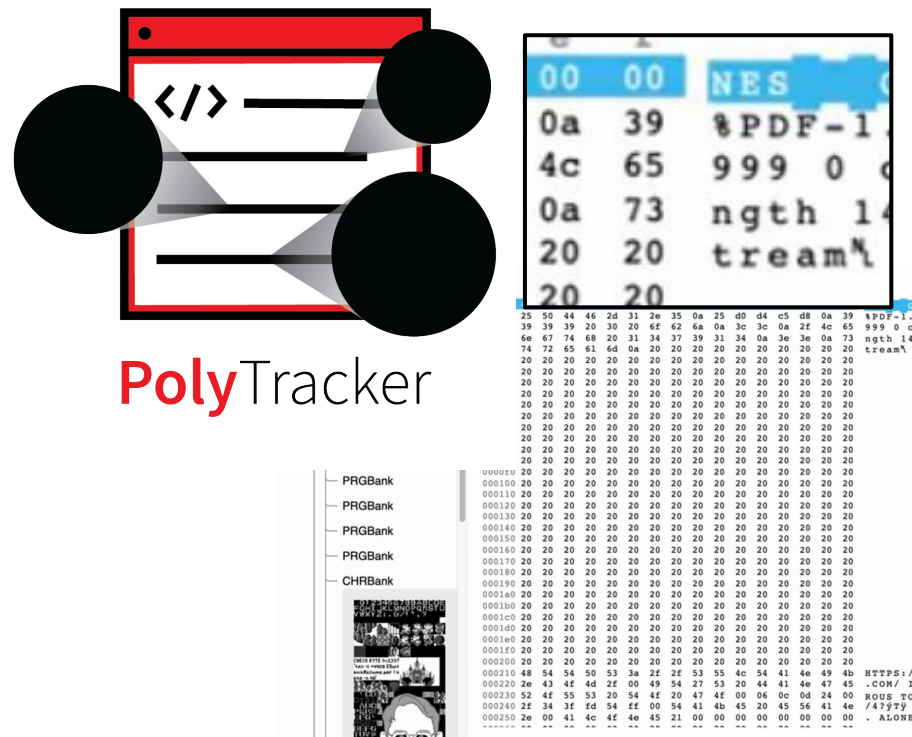Galois SafeDocs Structural Overview

# Identifying Sublanguages

- Large search space, unclear goals.

- Leverage known design principles.

- Three approaches surveyed:
  - **Taint Tracking**

  - **Entropy-based Methods**

  - **Reinforcement Learning**

# Taint Tracking

- Instrument existing parsers to link program logic and file bytes.

- Taint forest used to track complex interdependencies (non-CFG).

- Improved by adding ground truth (PolyFile, PolyMerge).

- Differential analysis.

- Understanding existing vs building unified, de facto parser.



**Poly**Tracker

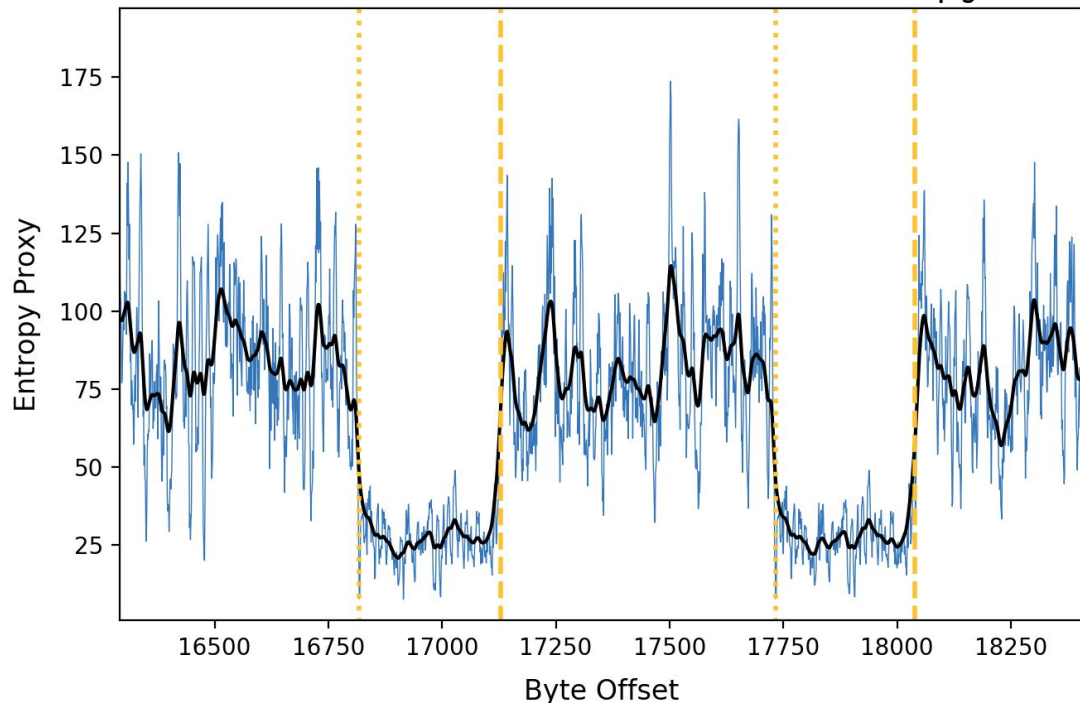# Entropy-Based Methods

- Parserless

- Data fingerprints

- Example result: PDF streams

Byte values

$$\frac{\sum_{i=0}^{n-2} |S_i - S_{i+1}|}{n-1}$$
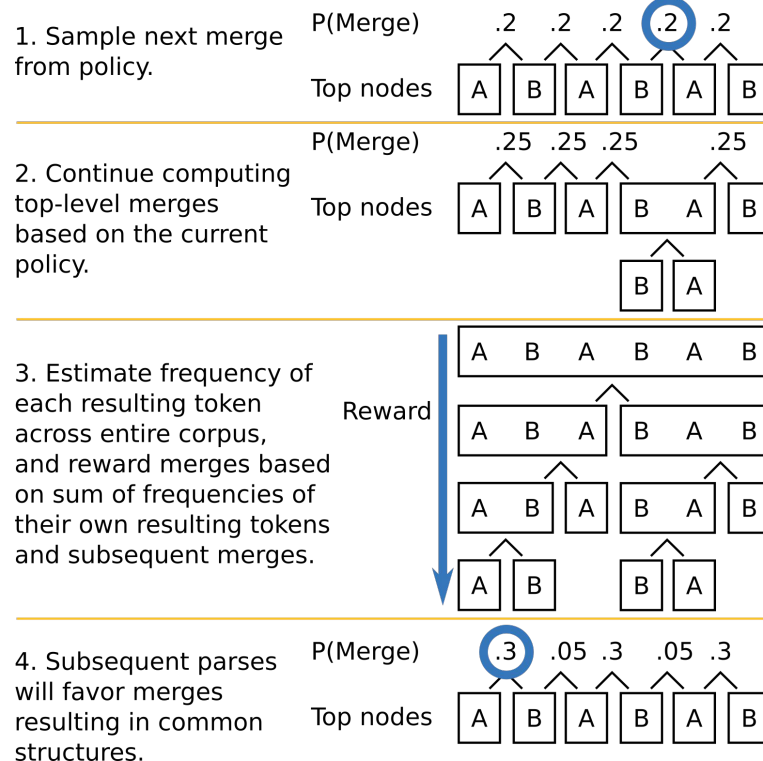
Window size

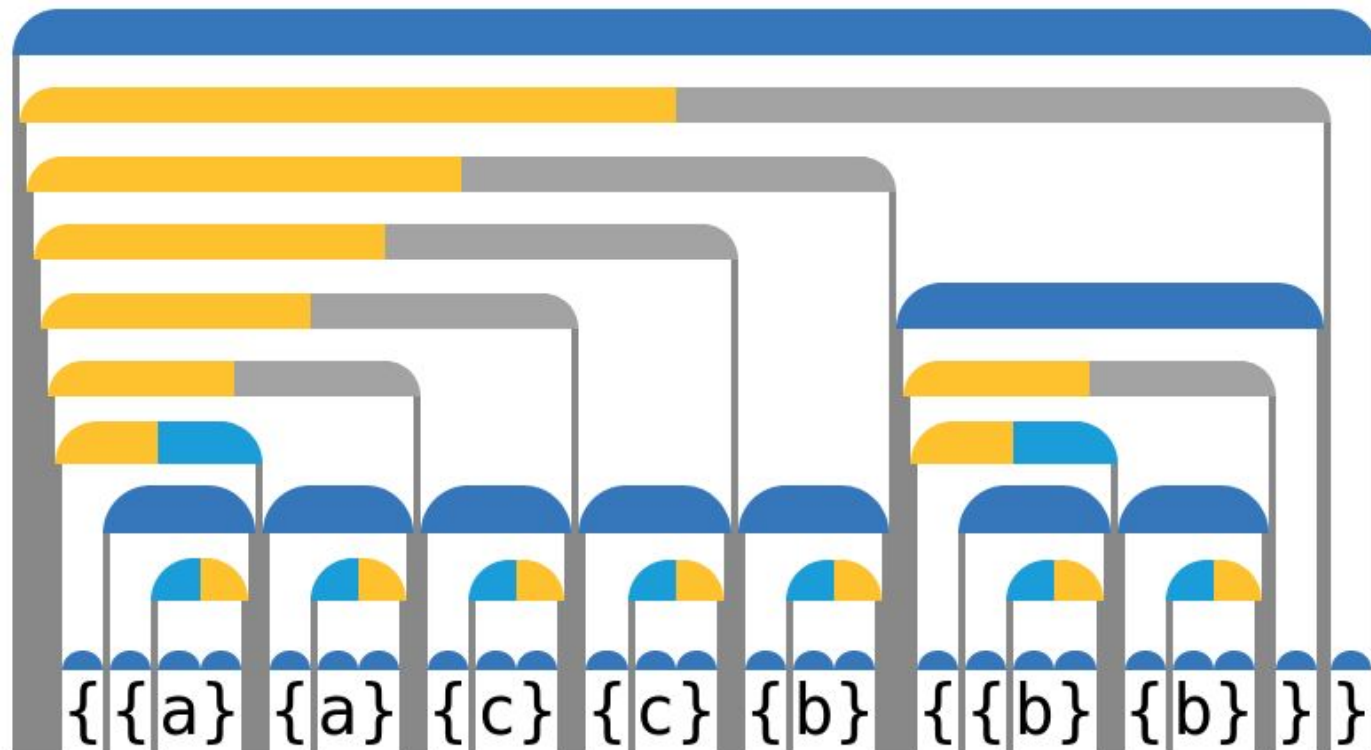### Position in File vs Measured Entropy

# Grammar Inference via Reinforcement Learning

- State-of-the-art grammar inference is focused on NLP, not data (outside of LearnPADS).

- Our research has focused on bottom-up parsing algorithm, using RL+statistics to derive parsers.

- RL provides flexibility for grammars outside of CFGs.

Learning to parse (AB)* via RL, example string ABABAB

1. Sample next merge from policy.

P(Merge)   .2   .2   .2   .2   .2

Top nodes   A  B  A  B  A  B

2. Continue computing top-level merges based on the current policy.

P(Merge)   .25  .25  .25       .25

Top nodes   A  B  A  B  A  B
                        B  A

3. Estimate frequency of each resulting token across entire corpus, and reward merges based on sum of frequencies of their own resulting tokens and subsequent merges.

Reward

A B A B A B

A B A  B A B

A B A  B A B

A B   B A

4. Subsequent parses will favor merges resulting in common structures.

P(Merge)   .3   .05  .3   .05  .3

Top nodes   A  B  A  B  A  B

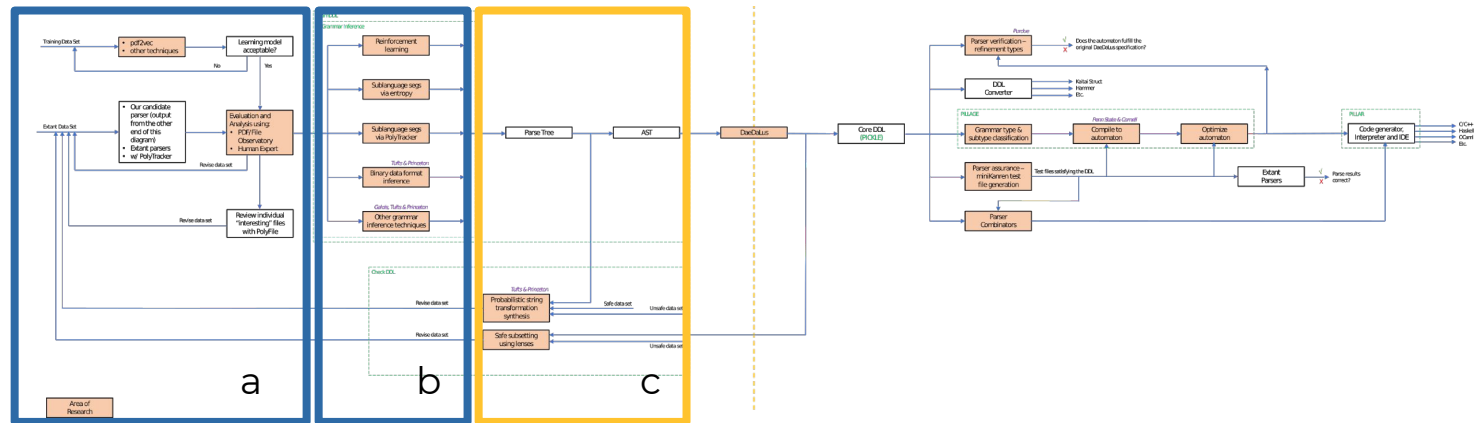# Grammar Inference via Reinforcement Learning



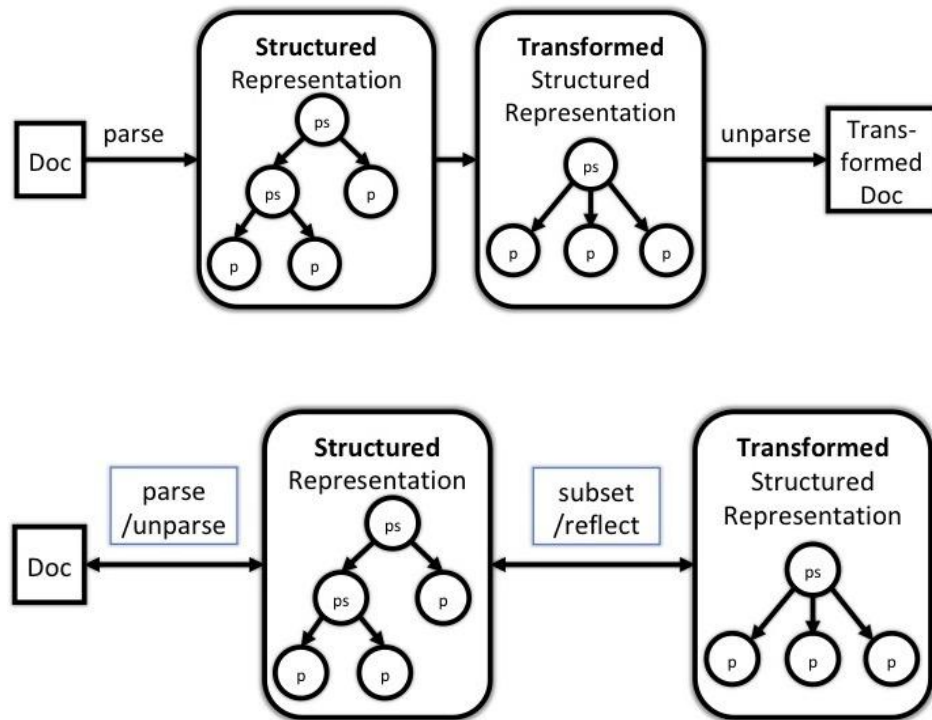How might we handle e.g. streams?

# ICARUS Project

- Pipeline and tools focused on developing methodologies for:
  a. Discovering and describing de facto data formats;
  b. Identifying grammars within a de facto format; and
  c. Translating from a de facto format to a safe subset.

Galois SafeDocs Structural Overview

# Safe Subsetting

- Desire to accept as much of de facto population as possible, while providing safety.

- Traditional subsetting.

- Bidirectional programming approach.

# Safe Subsetting

1. << /Size 14 /Version 5 >>

2. << /Size 14 /Version 5 /Version 6 >>

3. << /Size null /Version 5 >> == << /Version 5 >>

4. << /Size 14 /Version 5 /Version null >> == **??**

5. Lens ordering:
   - removeNullEntries . rejectDictionaryDups
   - rejectDictionaryDups . removeNullEntries

# Conclusions and Future Work

- ICARUS Toolchain being assembled to understand and secure de facto formats through:

  - Leveraging existing parsers,

  - Inferring de facto grammar,

  - Safe subsetting.

Portable Document Format Reference Manual

PDF Software

PDF

ICARUS