# ENGINEERING A XML-BASED CONTENT HUB FOR ENTERPRISE PUBLISHING

**Elias Weingärtner**

Christoph Ludwig

# HAUFE GROUP – QUICK FACTS

- Software Company and Media Publishing House
- Head Office: Freiburg, Germany
- Business Domains: Law, Tax, Human Resources, Talent Management, Trainings
- 150 Software Developers

# HAUFE: THE ROOTS

**Loose-leaf editions**

**Desktop content databases (1990s)**

**Books**

**Haufe.Gruppe**

**Online Content Databases**

**Haufe.de Portal Site**

**Booking platforms
for seminars & trainings**

**Books & Print Products**

❯ Weingaertner, Elias; Ludwig, Christoph - Engineering a XML-based Content Hub for Enterprise Publishing

# CONTENT @ HAUFE

**HAUFE.Gruppe**

- **50 million XML documents (Haufe Content)**

  - Own set of domain-specific DTDs

  - Proprietary Python-based publishing pipeline

    - Conversion to XML

    - Conversion to target formats (PDF, Database files)

- **Auxiliary content: PDFs, audio-visual content, forms, embedded applications**

- **News Posts**

- **Seminar descriptions**

# PROBLEM:
# SATURATED CONTENT MANAGEMENT

**iDesk2 App**

**Haufe Suite**

**Haufe.de**

Search
Retrieval

Search
Retrieval
Semantics

Search
Retrieval

Similar
Content

Similar Content

**Content Retrieval
System**

**L4**

**CoreMedia**

**Acquired Content
Bought-In Content**

**Content brokered
for other companies**

# PROBLEM:
# SATURATED CONTENT MANAGEMENT

1. **Complicated Content Reuse / Cross-Referencing**

2. **Difficult Authorization**

3. **Massive Content Duplication**

4. **High System heterogeneity → Increased management efforts**

# Vision: Unified Content Hub

# FUNCTIONAL BUILDING BLOCKS

**Indexing**

**Consistency**

**Content Storage**

**Map content structure to triple store → Integrity**

**Content graph for filtering / enhancing search**

**Search**

**Triple Store**

# CONTENT HUB ARCHITECTURE

Content Consuming Systems

...

Content Access Interface (CMIS)

Metadata Interface (SPARQL)

Search Interface & Query Processor

Authorization

Transformation

Aggregation

Transaction Management

Validation, Extraction & Transformation

Ingest Authorization

Single Document Ingest

Bulk Ingest

Content Sources

...

> Weingaertner, Elias; Ludwig, Christoph - Engineering a XML-based Content Hub for Enterprise Publishing

# WHY TRIPLES?

Products

Individual

Bundling

Construction Plan

Construction Plan

Construction Plan

Version 5

Most Recent

Construction Plan

Books

News

Seminars

Content

> Weingaertner, Elias; Ludwig, Christoph - Engineering a XML-based Content Hub for Enterprise Publishing

**Enables fast answers to complex questions**

- Display all seminars that discuss „Neuroleadership?"

- Enable cross references from free content (news posts) to relevant paid products

**RDF and triples for modeling relationships**

**SPARQL 1.1 for graph traversal**

**HAUFE.**Gruppe

**<link.norm bezeichner**="paragraph" **kuerzel**="EStG" **zahl**="32"**>**

*§ 32 des Einkommensteuergesetzes*

**</link.norm>**

**<link.text zielid**="HI39751.gen1"**>**

Über dieses Dokument

**</link.text>**

**<kuerzel basis**="Einkommensteuer-Richtlinien 1999"**>**

*Einkommensteuer-Richtlinien 1999*

**</kuerzel>**

# IMPLEMENTATION OPTIONS

**Use Document Repository**

Think Fedora or JCR

**Scale-Out?**

**XML-agnostic**

**→ No XPath, XQuery…**

**Integrate XML-DB**

**With scalable
search & triple store**

**Certainly doable**

**Core Problem: Data Integrity**

**Use XML-based,
integrated solution**

**Don't reinvent the wheel**

**Focus on developing
add-on functionality**

**License fees**

**Dependency on vendor**

> Weingaertner, Elias; Ludwig, Christoph - Engineering a XML-based Content Hub for Enterprise Publishing

# TIMELINE: PAST, PRESENT, FUTURE

**September 2013:**     **Business department wants TWO new systems:**

- Global Content Search
- Unified Content Hub

**Fall 2013**     **Three Software architects create two architectural drafts**

**Outcome:**     Search without docs? Store without search?

Data Integrity? How to deal with graph structure?

**Winter 2013/2014**     **Consolidation of Drafts → Unified Content Hub**

**Spring 2014**     **Proof of Concept with major XML NoSQL vendor**

- Identification of additionally required external services
- Further elaboration of triple use

**Summer 2014-**     **Start of Implementation**

# SUMMARY & CONCLUSION

**Haufe.**Gruppe

- **Consolidation of saturated storage and search services**

  → Avoid content duplication

  → No duplicated indexing

  → Reduce infrastructure and management costs

- **Indexing XML Structure is vital**

  → Faceted search & complex search using XPath / XQuery

- **Triples for relationship management**

  → Will allow querying structure in real-time

  → Triples for modeling

  → SPARQL1.1 for querying and graph traversal

- **Currently working towards first implementation**

> Weingaertner, Elias; Ludwig, Christoph - Engineering a XML-based Content Hub for Enterprise Publishing