

## A Quantitative Comparison of Binary XML Encodings

Wolfgang Hoschek

Lawrence Berkeley National Laboratory

[whoschek@lbl.gov](mailto:whoschek@lbl.gov)

- XML established for interop. data exchange
- Grid and Web services
  - Large numbers of XML messages at high frequency
  - XML serialization & deserialization bottlenecks
    - Markup verbose in size
    - Expensive to produce and parse
    - Restricts XML adoption
  - Alternative: Binary XML?
- Novel binary XML encoding (*bnux*)
- Quantitative evaluation
  - Production-quality XML and Binary XML toolkits
  - Tree and streaming deserialization, serialization, compression

- **Faithful to XML**
  - General purpose
  - Preserves all information without loss or change
  - Preserves W3C XML InfoSet and W3C Canonical XML
- **Self-contained**
  - No external resources required (e.g. no schema)
- **Tree and streaming deserialization mode**
  - For end user applications and filter pipelines
- **Tunable for either performance or size**
  - High vs. low bandwidth networks
  - Moderate compression via simple FAST means (tokenization)
  - Additional strong GZIP (ZLIB) compression (optional)
- **Production quality implementation**
  - Serialization of XOM XML object model (~DOM)
  - <http://dsd.lbl.gov/nux>

## - Serialization

Extract unique symbols (strings) via hash table

Sort symbols by frequency (top N)

Encode symbol table as zero terminated UTF8

Encode each XML node as binary token, with compact pointers into symbol table (Vint)

## - Deserialization

Decode UTF-8 symbol table

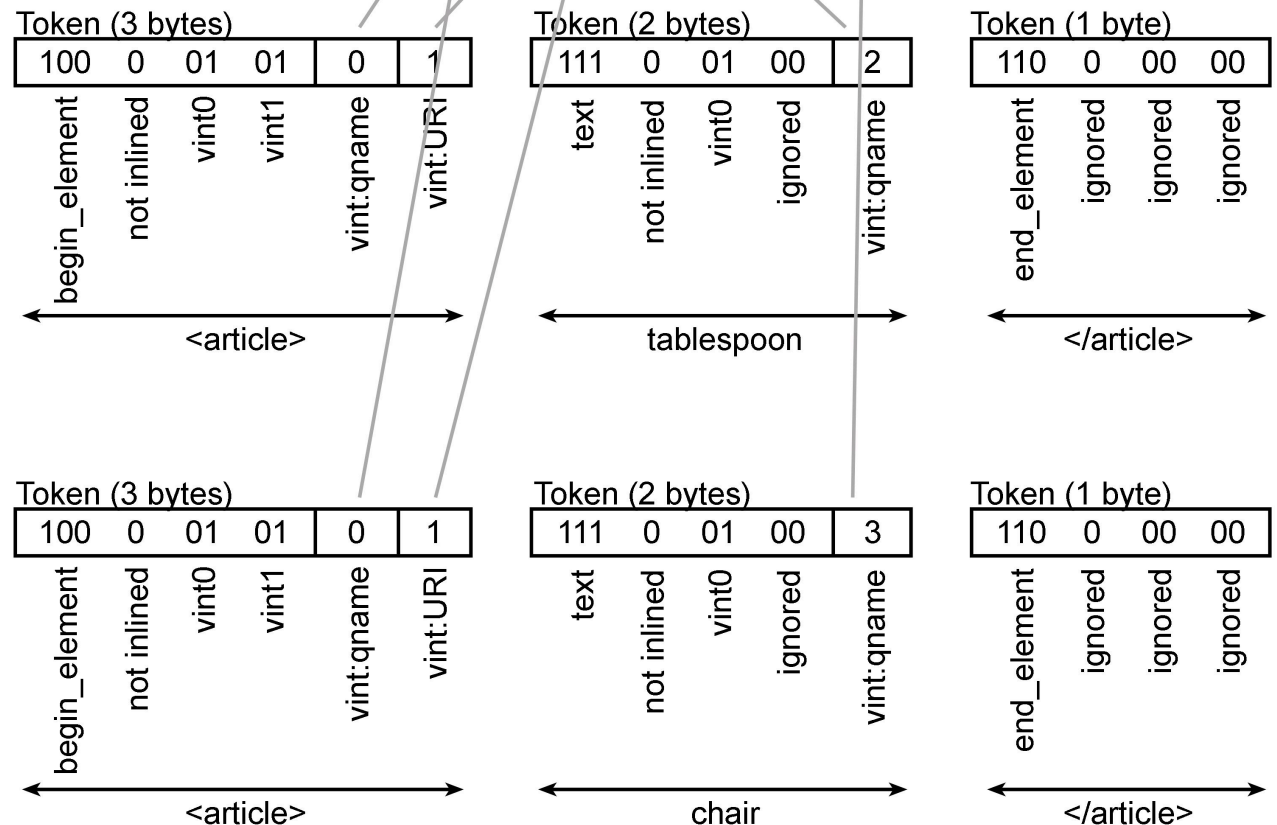
Decode tokens, hand info to app handler (e.g. type, prefix, name, URI)

Example XML (53 bytes)

```
<article>tablespoon</article><article>chair</article>
```

Serialized Symbol Table (26 bytes)

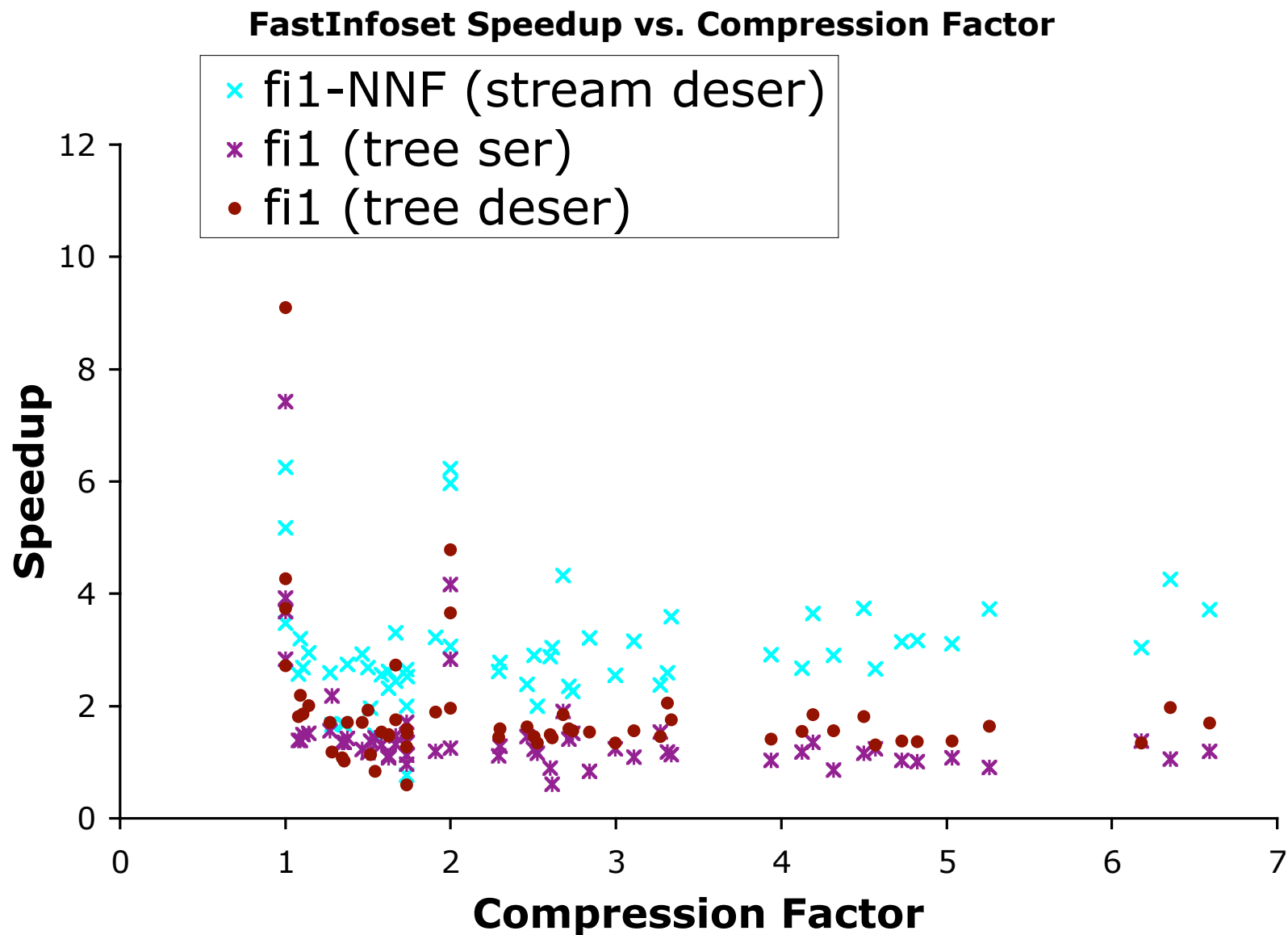
```
"article"0, ""0, "tablespoon"0, "chair"0
```



- **Workloads: 60 distinct test document flavours**
  - Wide range of real-world documents
  - File size
    - Small (0.2 - 1 KB), medium (1KB - 4 MB), large (4 - 100 MB)
  - Documents
    - Messaging-oriented, record-oriented (database), narrative text
    - E.g. WSDL, SOAP, RSS, ATOM, DB, Shakespeare, P2PIO, ...
  - With and without namespaces, attributes, whitespace, repetitions, nesting depth, ...
- **Memory-to-memory tests (no I/O perturbation)**
- **Setup**
  - Sun Java 1.5.0\_04, server VM, PentiumIV Xeon@2.8 Ghz, 2GB memory, Linux 2.4.20
  - xom-1.1, nux-1.4, saxonb-8.5.1, java.net-fastinfoset-CVS (ISO/ITU), xerces-2.7.1 for SAX and DOM, woodstox-2.0.2 for STAX

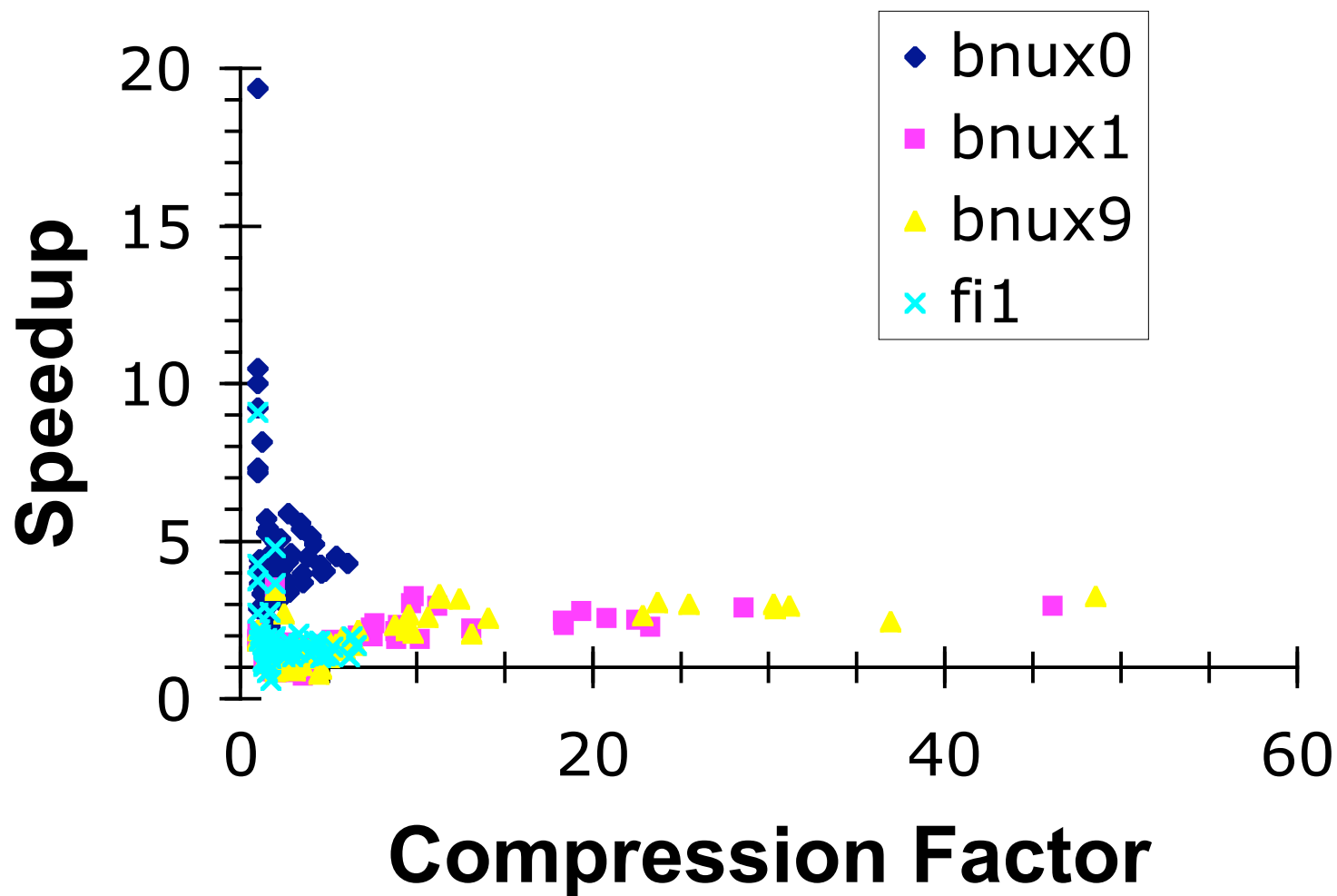
Model	Description
xom-NV	XOM via SAX/Xerces with XML verification performed by Xerces. <b>Comparison baseline</b> for tree speedup.
xom-V	Same as xom-NV except with XML verification performed by XOM instead of Xerces. More expensive than xom-NV
saxon	Saxon tinytree model with shared namepool (via SAX/Xerces)
dom	Xerces Document Object Model without “deferred node expansion”
bnux0	Bnux binary XML with XML verification; no GZIP compression
bnux0-NV	Same as bnux0, except that PCDATA verification is omitted
bnux1	Same as bnux0, plus weak GZIP compression at level 1
bnux9	Same as bnux0, plus strong GZIP compression at level 9
fi0	FastInfoset binary XML with default indexing (via SAX)
fi1	FastInfoset binary XML with “full indexing” feature (via SAX)
xom-NNF	<b>Streaming</b> XOM via SAX/Xerces with NullNodeFactory handler, throwing away all data, building an empty tree instead. <b>Comparison baseline</b> for streaming speedup.
bnux0-NNF	Same as xom-NNF except that bnux0 is used; no verification
fi0-NNF	Same as xom-NNF except that fastinfoset is used



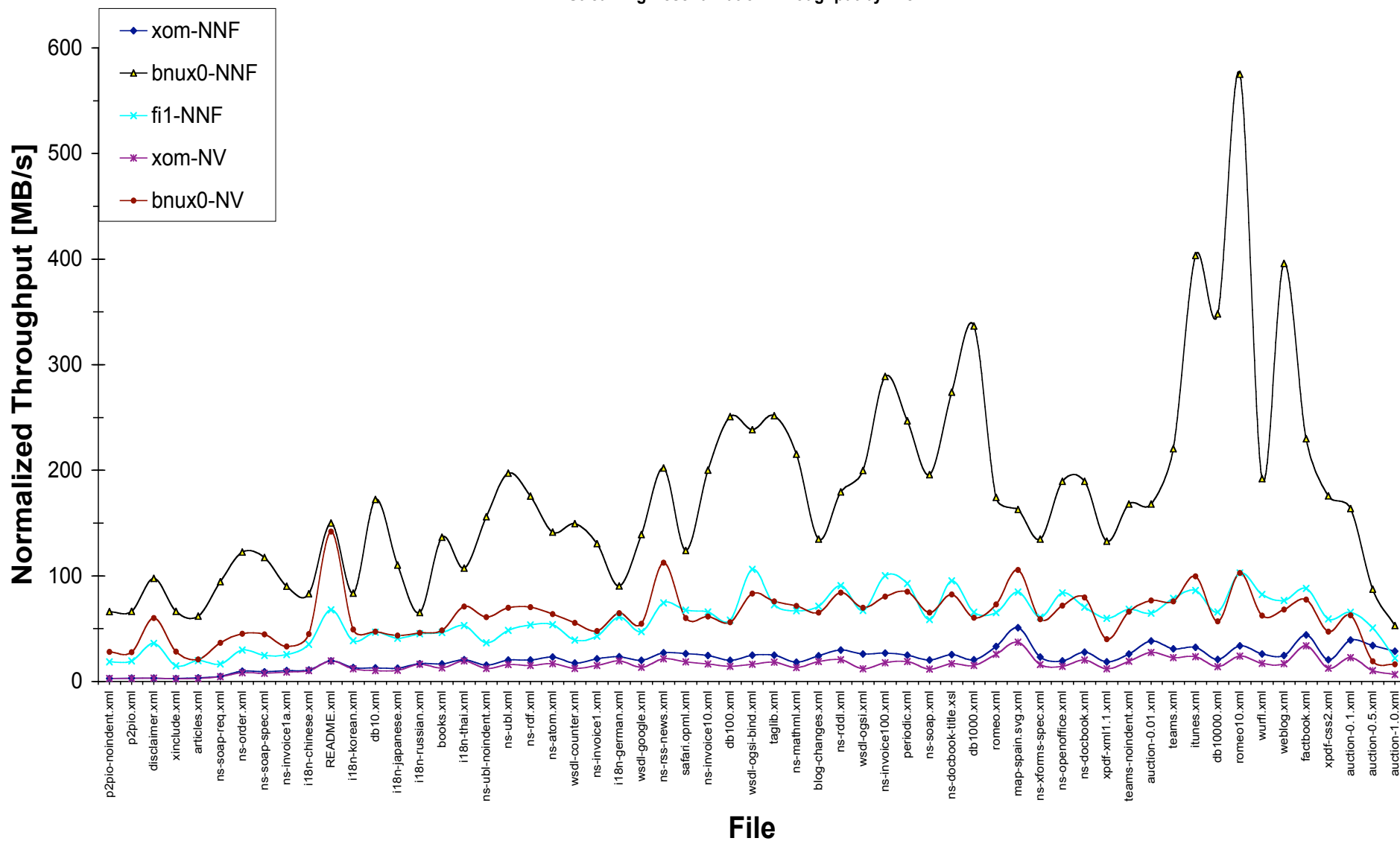




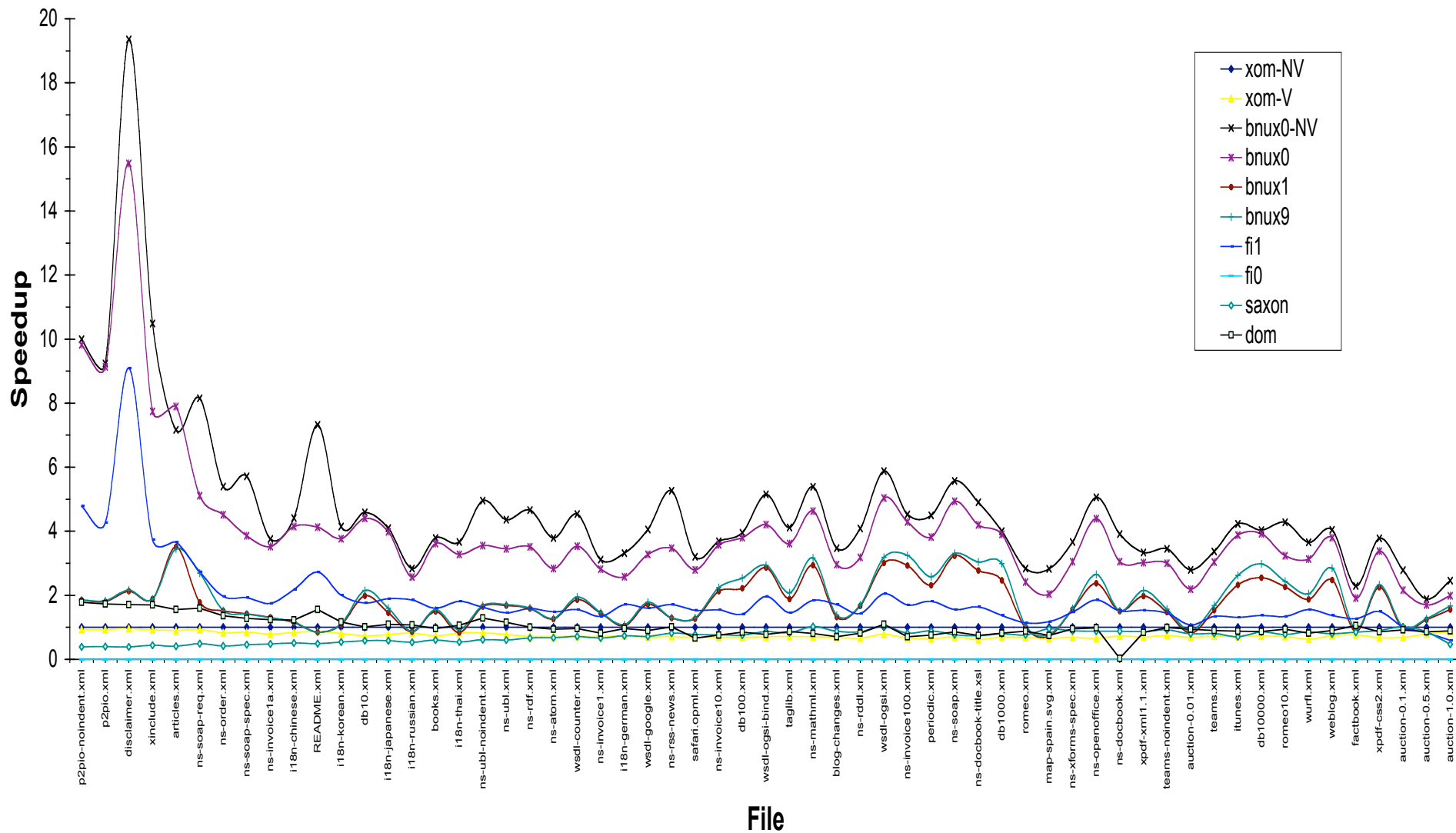
Tree Deserialization Speedup vs. Compression Factor



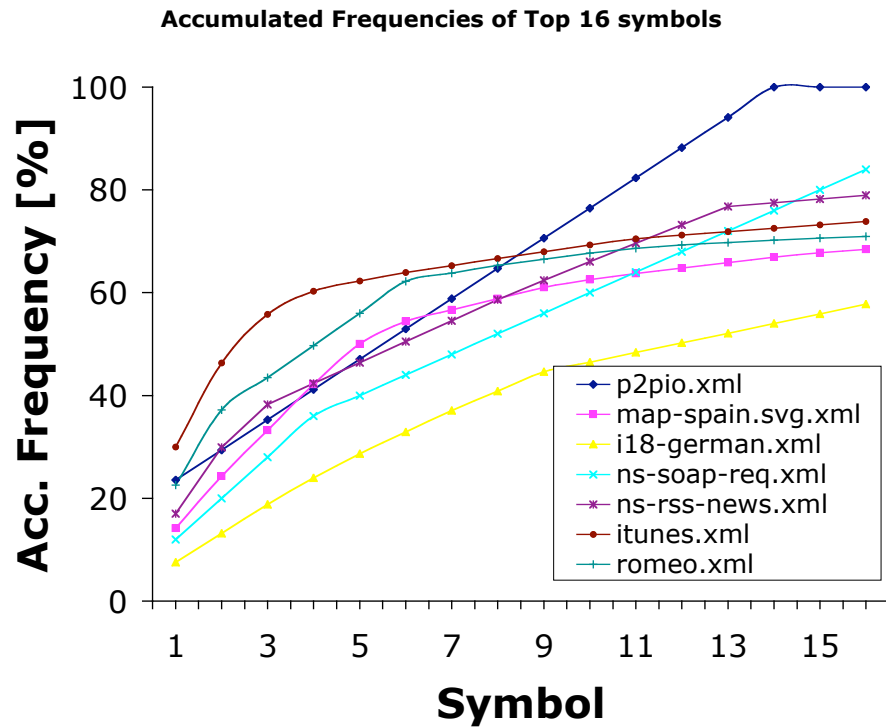
Streaming Deserialization Throughput by File



Tree Deserialization Speedup by File

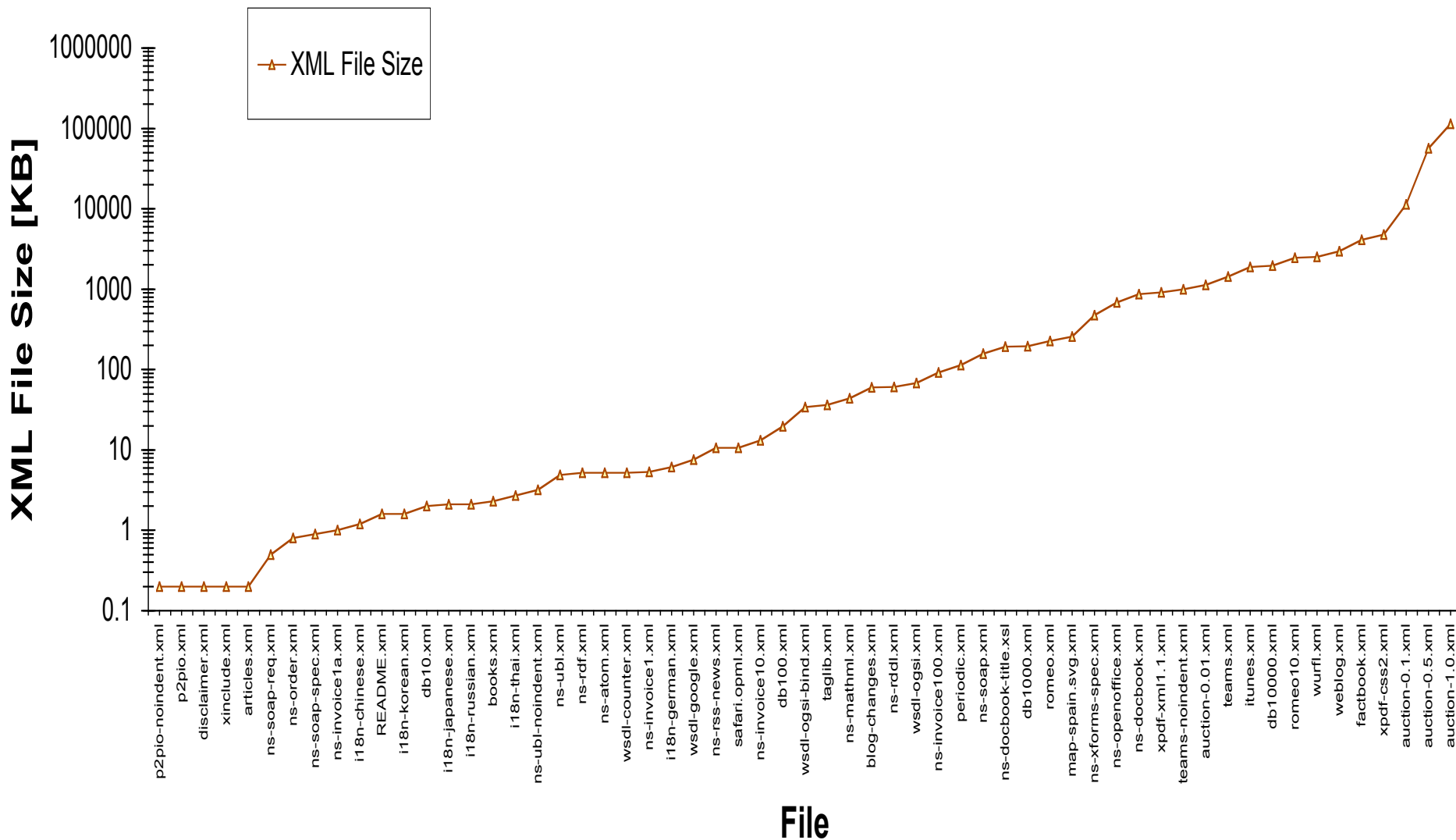


- **Standard textual XML satisfactory for many commodity use cases**
  - But requires non-intuitive configuration wizardry
- **Binary XML significantly faster for demanding data-intensive use cases**
  - Small to medium sized messages: 5-20x
  - Large variance stemming from document flavour
- **Moderate compression via tokenization (1-5x) is fast**
- **Strong compression via GZIP (5-50x) is too slow**
- **Trading efficiency for standardization?**
- **Input for potential W3C standardization?**

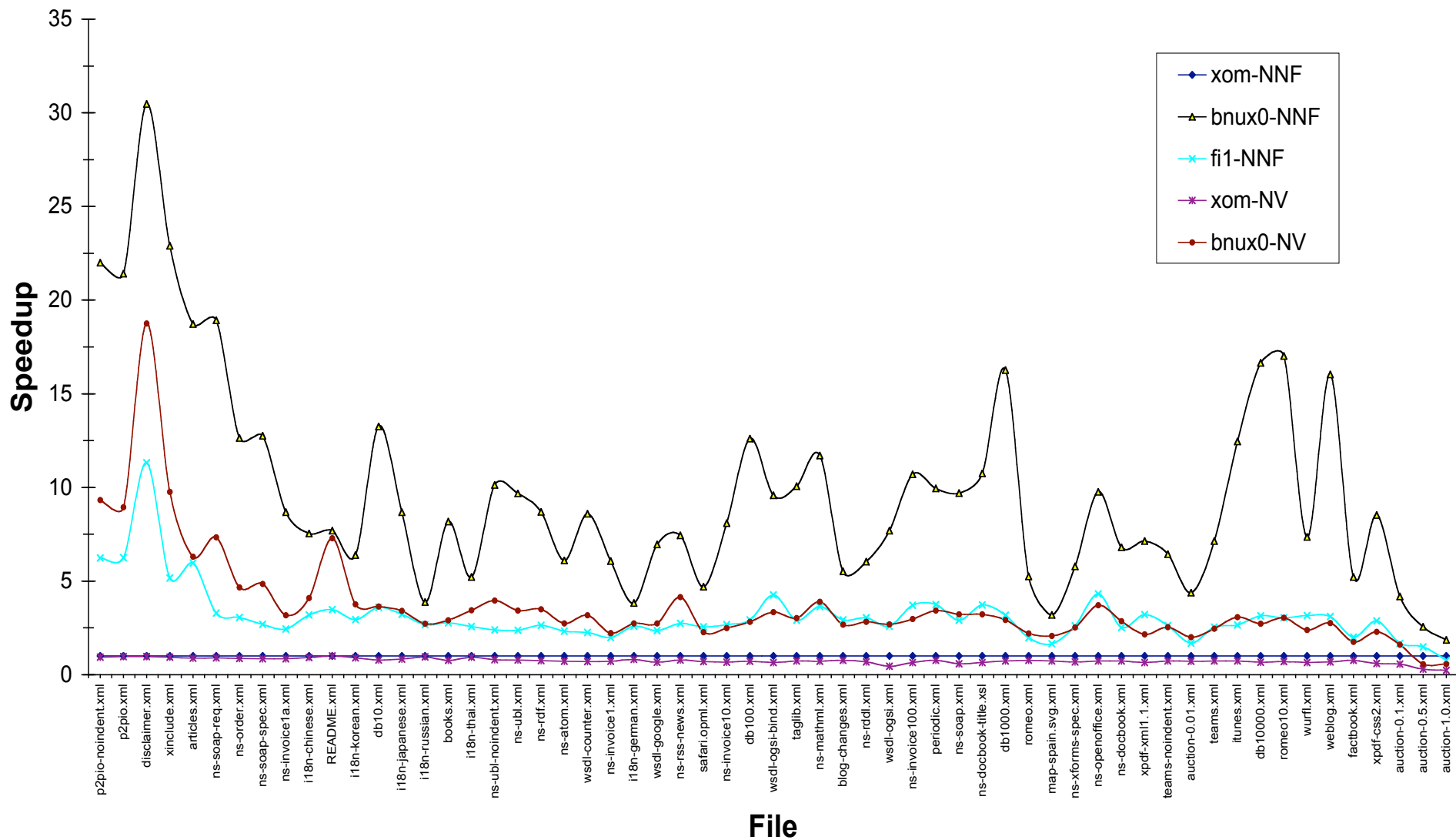


File	Compression Factor	Streaming Deser Speedup	XML File Size [KB]	Unique Symbols [%]
p2pio.xml	1	21.4	0.2	92.6
map-spain.svg.xml	1.2	3.2	258.4	81.5
i18n-german.xml	1.2	3.8	6.1	64.1
ns-soap-req.xml	1.3	18.9	0.5	80
ns-rss-news.xml	1.5	7.4	10.6	57.3
romeo.xml	1.5	5.2	228.2	97
itunes.xml	4.5	12.4	1882	83.7
romeo10.xml	6	17	2470	0.1

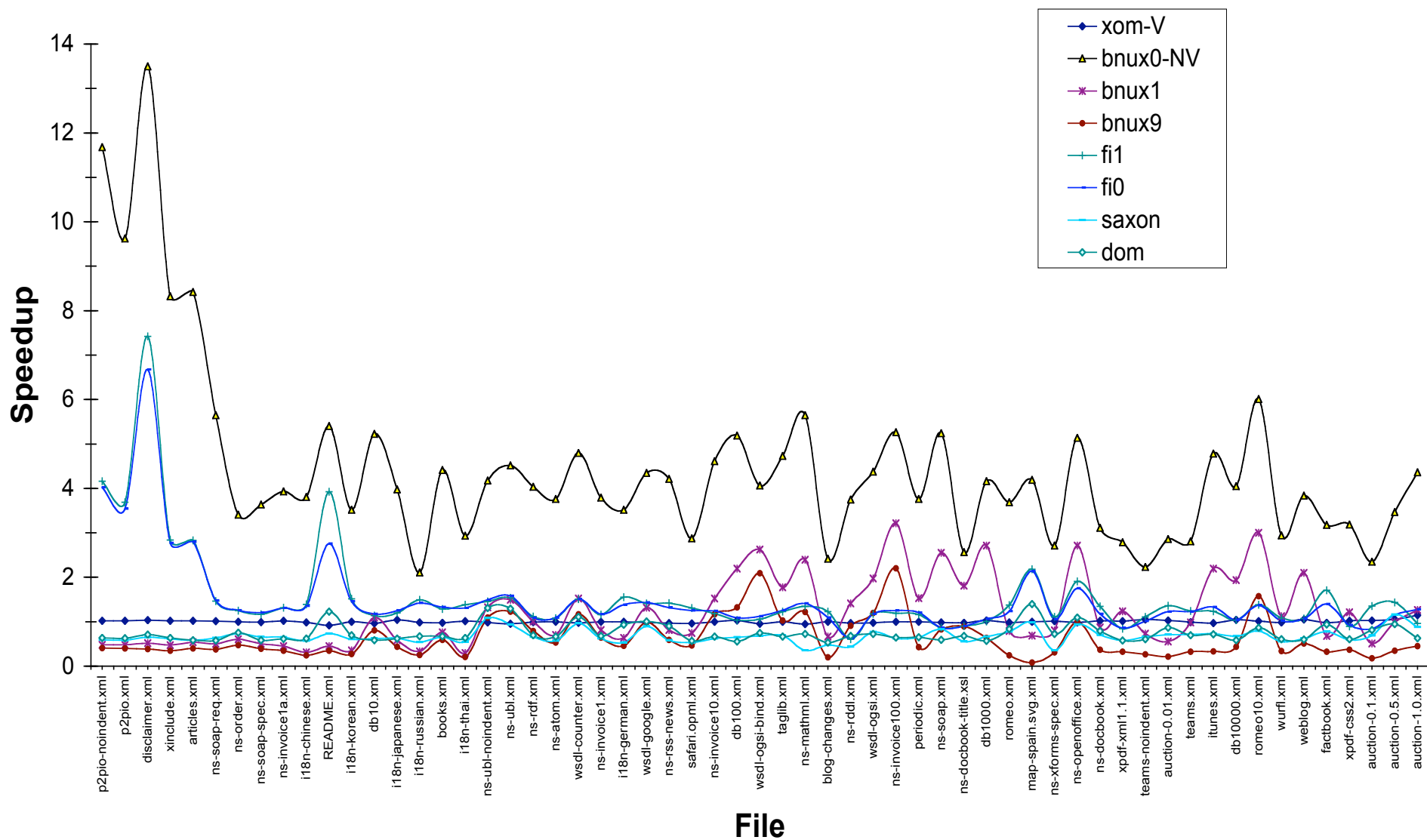
XML File Size of Document Flavours



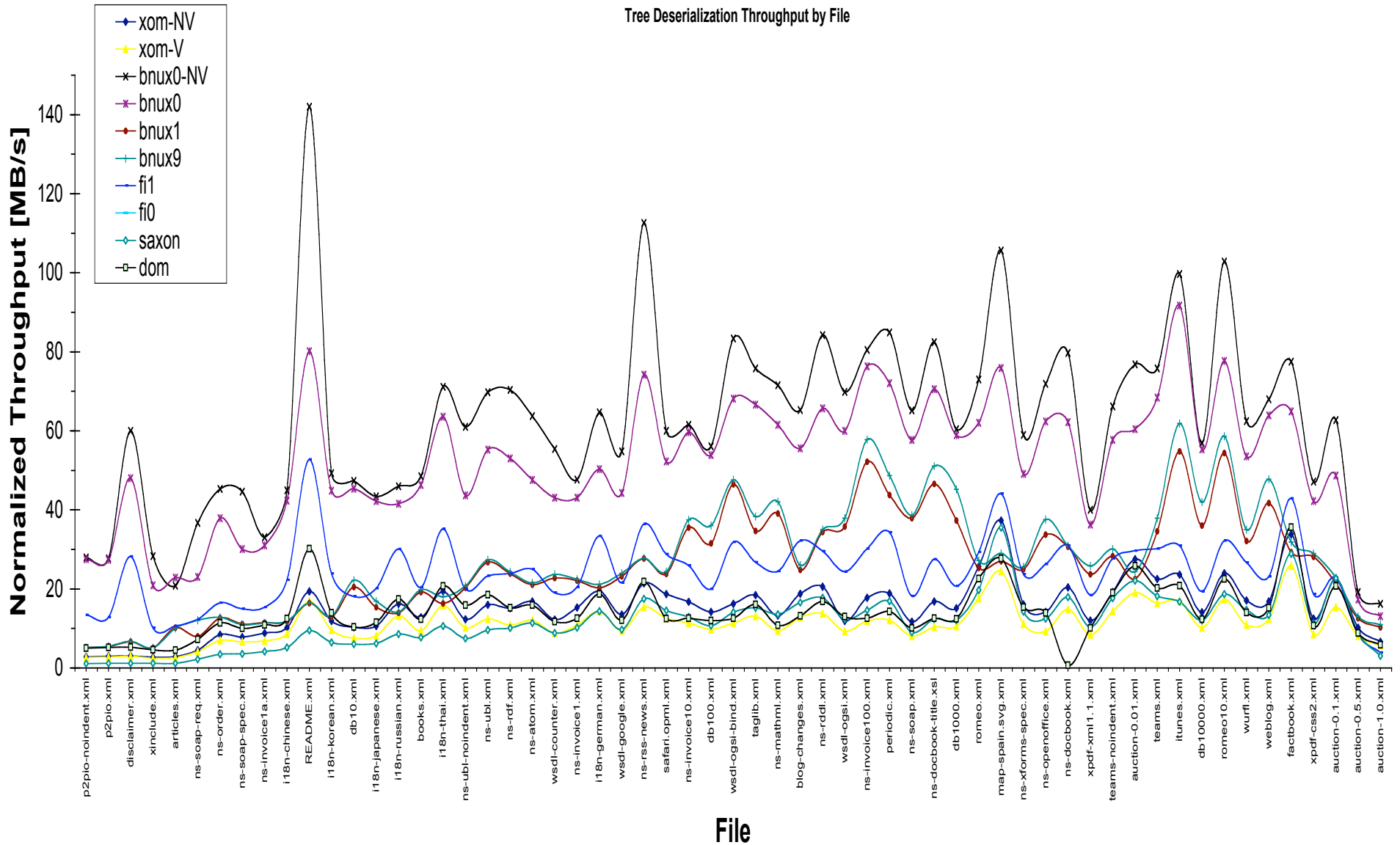
Streaming Deserialization Speedup by File

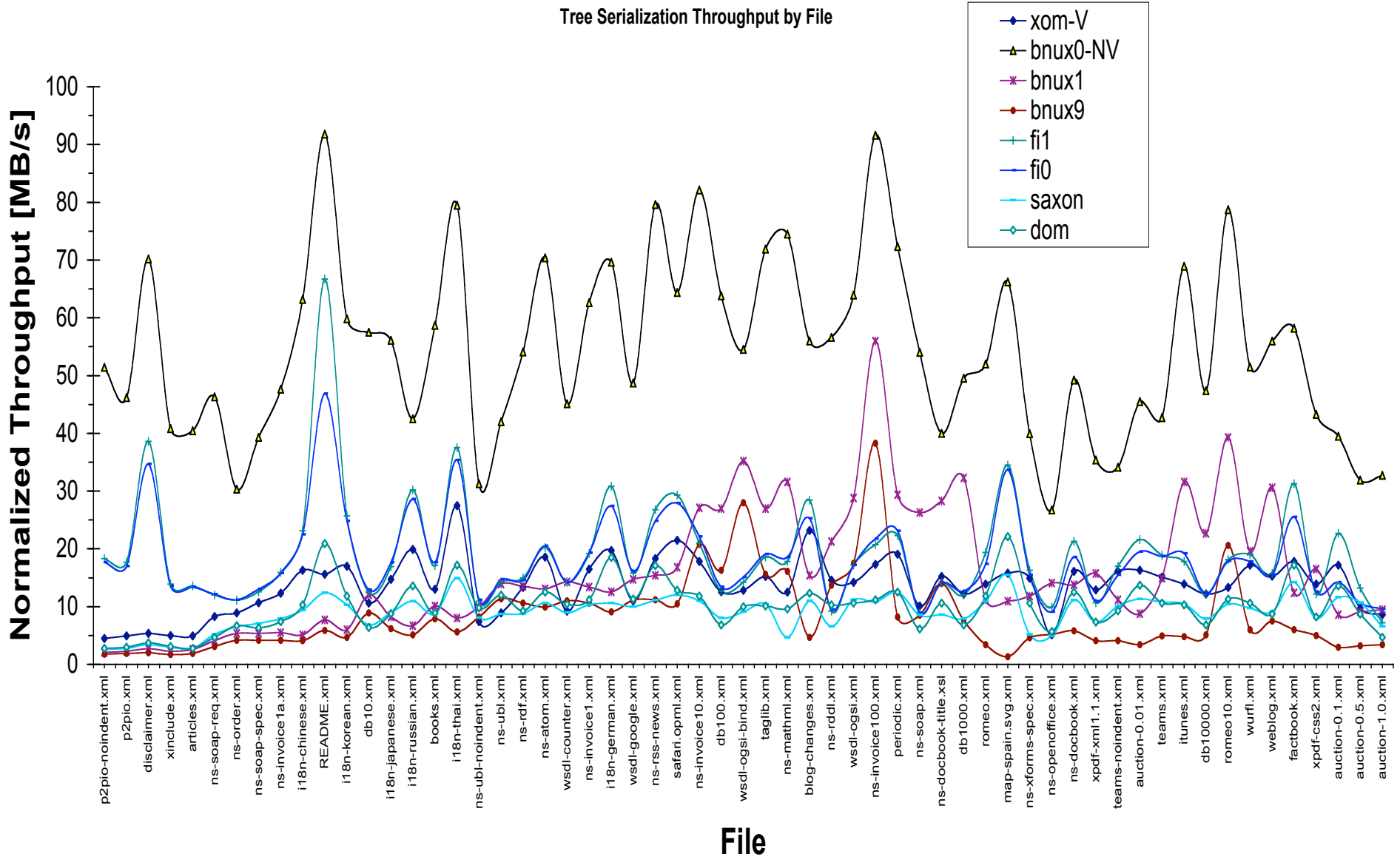


Tree Serialization Speedup by File

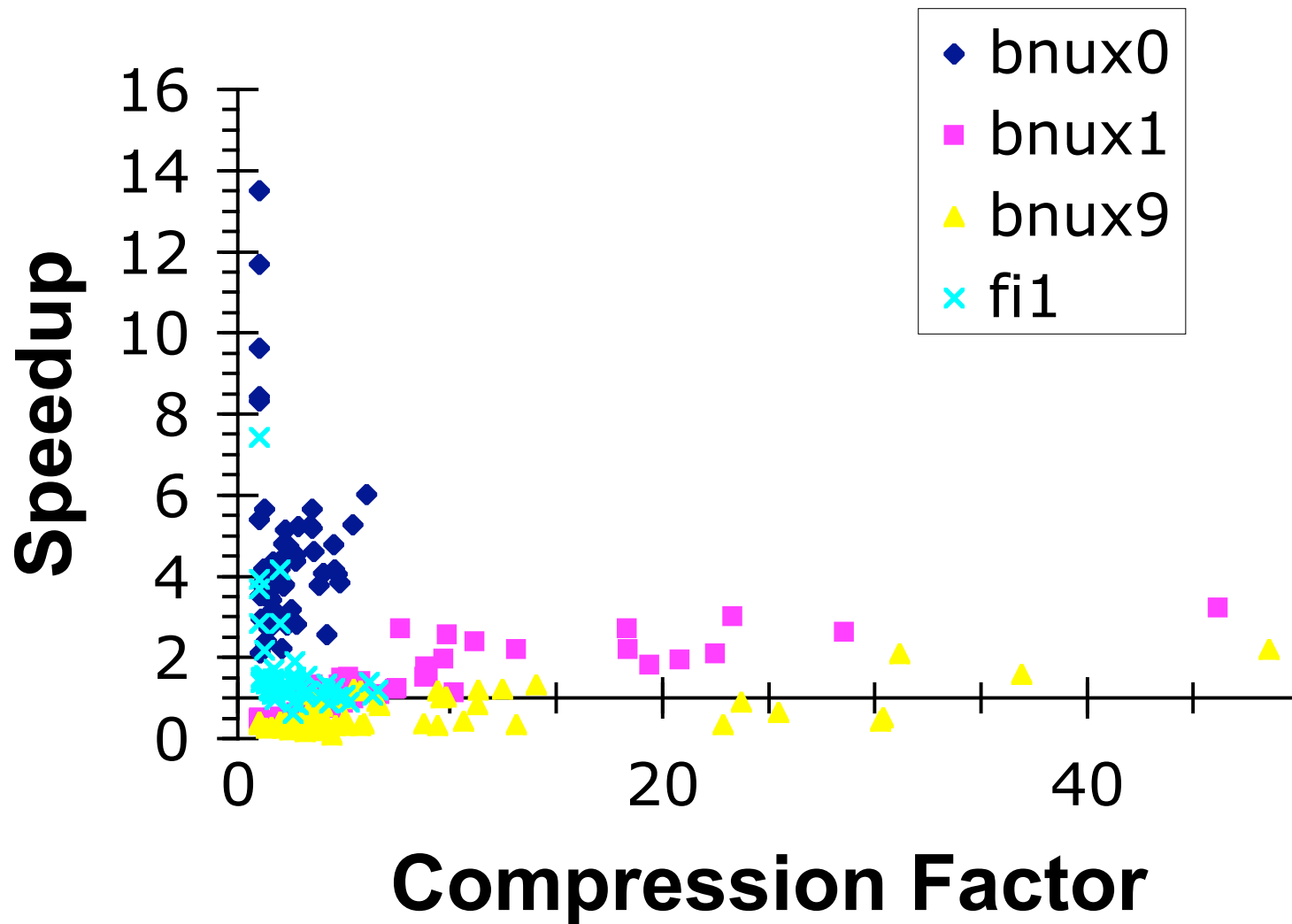




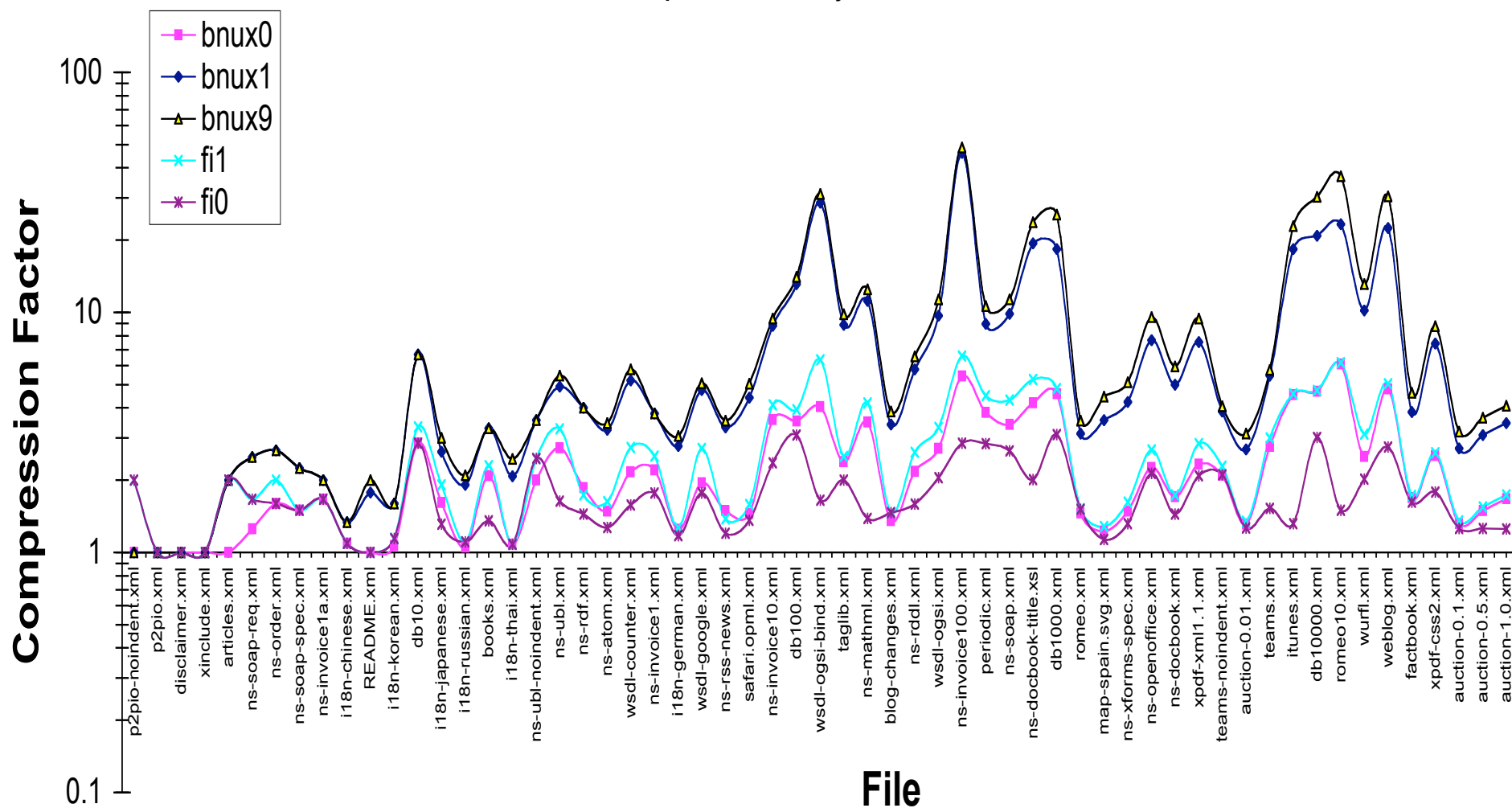




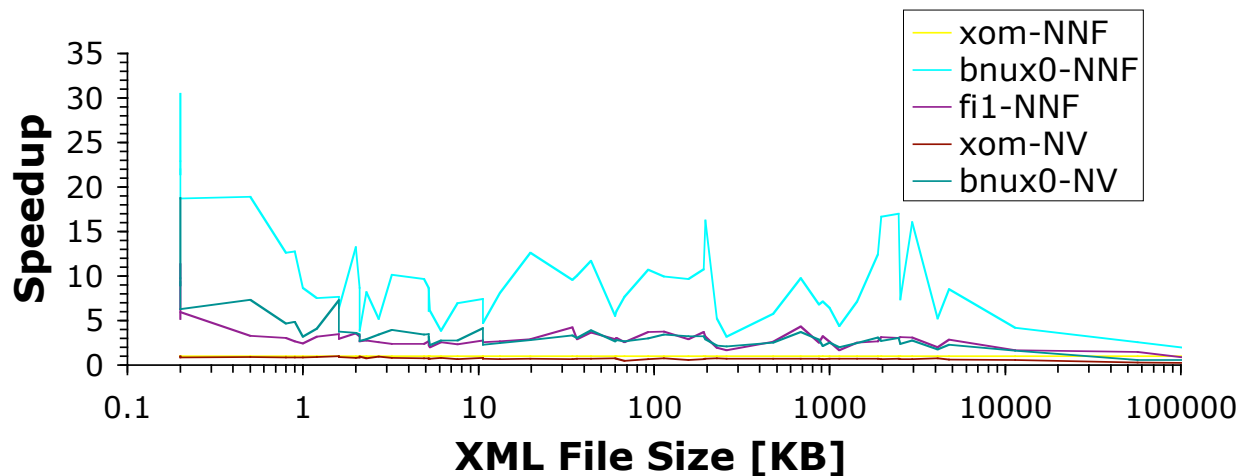
Tree Serialization Speedup vs. Compression Factor



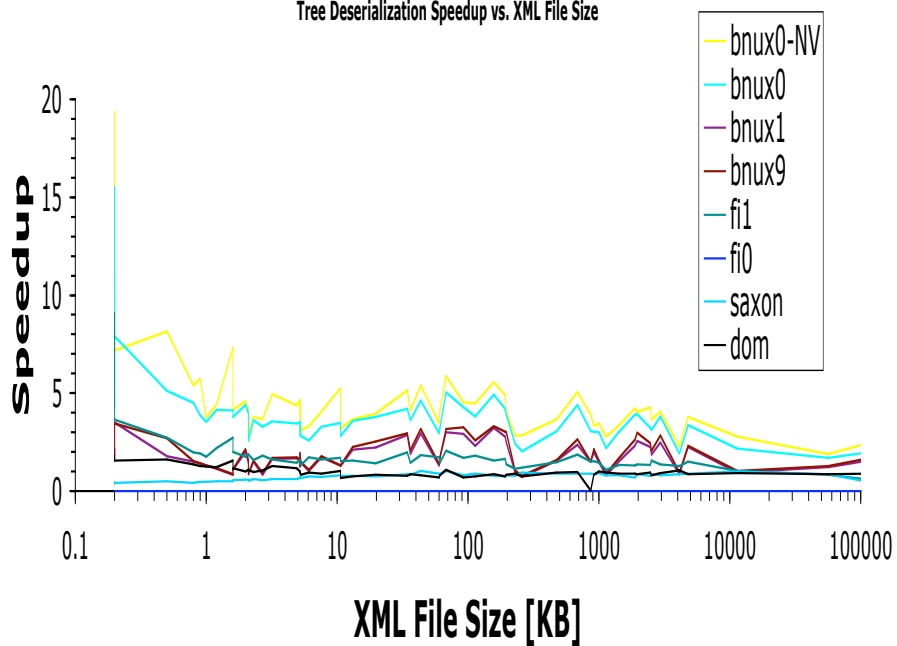
Compression Factor by File



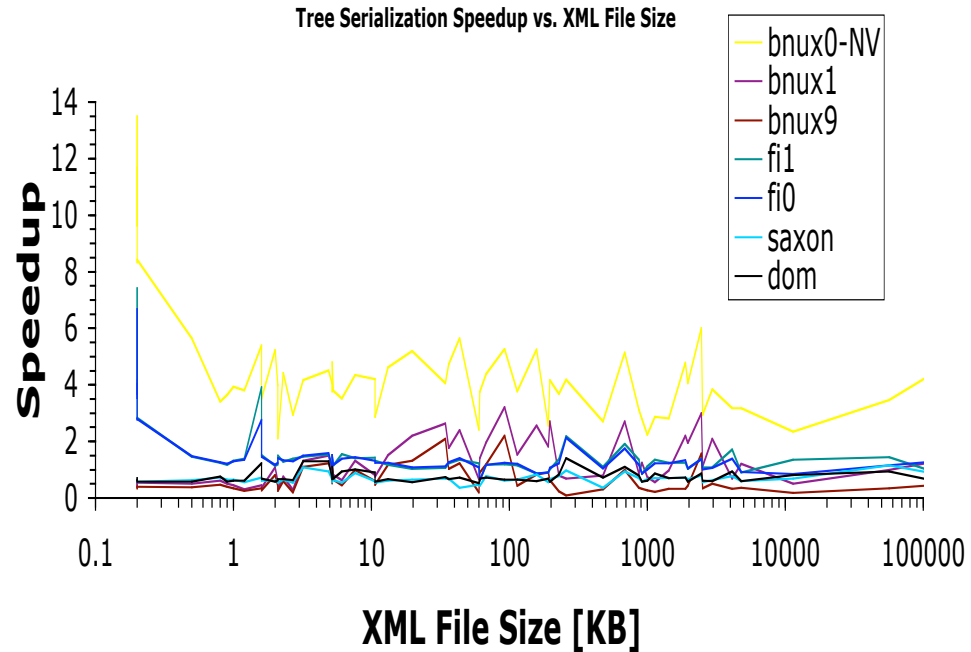
Streaming Deserialization Speedup vs. XML File Size



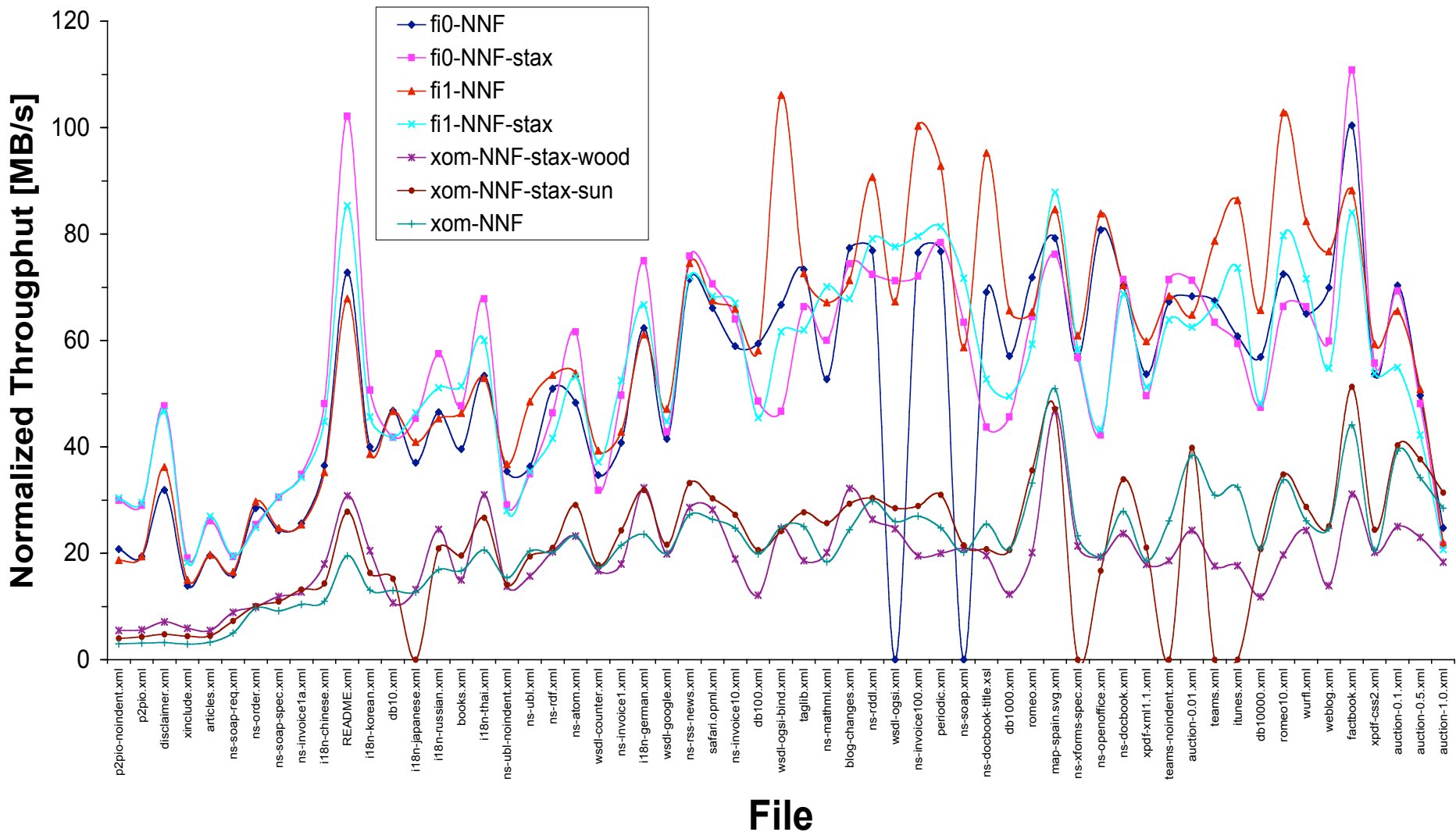
Tree Deserialization Speedup vs. XML File Size



Tree Serialization Speedup vs. XML File Size



## STAX vs. SAX Streaming Deserialization Throughput for FastInfoset & XOM



## STAX vs. SAX Tree Deserialization Throughput for FastInfoset & XOM

