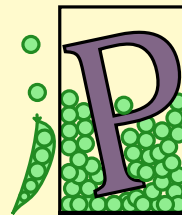


[title}

Markup Beyond XML

{title [alt}A LMNL Progress Report{]}



Wendell Piez

www.wendellpiez.com

Digital Humanities 2013

Lincoln, Nebraska

July 19 2013

Eat Your Vegetables

```
[lg]
[q [who]Maddalo{who}][l [n]96{n}]Look, Julian, on the west,
and listen well{1}
[lg]
[1 [n]97{n}]If you hear not a deep and heavy bell.{1}{q}
{lg}
[lg]
[1 [n]98{n}]I looked, and saw between us and the sun{1}
[1 [n]99{n}]On an island; such a one{1}
{lg}
[lg]
[1 [n]100{n}]As age, age might add, for seas {1}
[1 [n]101{n}]A windowless, deformed and dreary pile;{1}
{lg}
[lg]
[1 [n]102{n}]And on the top an open tower, where hung{1}
[1 [n]103{n}]A bell, which in the radiant vaults
[1 [n]104{n}]We could just hear its hoarse and iron tongue:{1}
{lg}
[lg]
[1 [n]105{n}]The bright sun sunk behind it, and it tolled{1}
[1 [n]105{n}]In silver and black relief.— [q [who]Maddalo{who}]What we behold{1}
{lg}
[lg]
[1 [n]107{n}]Shall be the madhouse and its belfry tower,{1}{q}
[1 [n]108{n}]Said Maddalo, [q [who]Maddalo{who}]and ever at
this hour{1}
{lg}
[lg]
[1 [n]109{n}]Those who may cross the water, hear that bell{1}
[1 [n]110{n}]Which calls the maniacs, each one from his cell,{1}
{lg}
[lg]
[1 [n]111{n}]For ever, [q [who]Maddalo{who}]and ever shall we read {1}
```

(pronounced “liminal”)

What is LMNL?

A data model

An approach to markup

An experiment in text encoding*

* An example of what can be done on a budget

Jeni Tennison and Wendell Piez. “The Layered Markup and Annotation Language (LMNL).”
Extreme Markup Languages 2002 (Montréal, Canada: August 2002).

The Caterpillar and Alice looked at each other in silence: at last the Caterpillar took a coil of itself round its middle and spoke in a languid, sleepy voice.

The LMNL Model

A document consists of **text**, a sequence of **atoms**
(Most atoms will be represented as characters)

With **ranges** (subsequences) over the text

Atoms and ranges may be named and annotated

Annotations may also be named and annotated

Text in annotations may have ranges

(Annotations are like documents)

Ranges over a text may have any (or no) relation
to one another

"It isn't," said the Caterpillar.

"Well, perhaps you haven't found it so yet," said Alice.

Implementation

Implicit

XML may be mapped directly to LMNL

Elements cover ranges

(or segments of ranges, or “milestone” markers at range start/end points)

Or, LMNL may be represented using standoff markup

Elements* maintained out of line indicate and represent ranges

Pointers indicate range locations and extent

Range or standoff properties (elements, attributes) become annotations

* Or some other representation or notation (yes JSON)

Out of line

Or, LMNL may be maintained in a database or object structure

represented in a UI but never serialized

Object model

But — range models have an Achilles heel:

Defining, deploying, controlling and maintaining the pointers ...

Plus, we will soon want a serialization (external format)

... well ...

(Markup syntax?)

Markup

... LMNL has a syntax ...

```
[poem [title]Paradise Lost{] [author]John Milton{]}
[book [n]1{n}]
[verse-paragraph]
[1 [n]1.1{]}[s]{phr}Of Man's first disobedience,{phr} [phr]and the fruit{1]
[1 [n]1.2{]}Of that forbidden tree whose mortal taste{1]
[1 [n]1.3{]}Brought death into the World,{phr} [phr]and all our woe,{phr}{1]
[1 [n]1.4{]}[phr]With loss of Eden,{phr} [phr]till one greater Man{1]
[1 [n]1.5{]}Restore us,{phr} [phr]and regain the blissful seat,{phr}{1]
[1 [n]1.6{]}[phr]Sing,{phr} [phr]Heavenly Muse,{phr} [phr]that,{phr} [phr]on the secret top
[1 [n]1.7{]}Of Oreb,{phr} [phr]or of Sinai,{phr} [phr]didst inspire{1]
[1 [n]1.8{]}That Shepherd who first taught the chosen seed{1]
[1 [n]1.9{]}In the beginning how the heavens and earth{1]
[1 [n]1.10{]}Rose out of Chaos:{phr} [phr]or,{phr} [phr]if Sion hill{1]
[1 [n]1.11{]}Delight thee more with Helicon's sweet brook that flowed{1]
[1 [n]1.12{]}Fast by the oracle to see the Pythian priestess dancing{1]
[1 [n]1.13{]}Invoke thy aid to me inspir'd, thy quickning powers impart{1]
[1 [n]1.14{]}[phr]That with my slumbers thou may'st wake, and with thy haughty
[1 [n]1.15{]}Above th' Aonian top, thy voice may'st flourish o'er thy art{1]
[1 [n]1.16{]}Things unattempt'd yet in prose or rhyme {phr}{s}{1]
[1 [n]1.17{]}[s]{phr}And chiefly thou, my spirit, thy presence gladdens me{1]
[1 [n]1.18{]}Before all temples thou shalt see, my soul to thee doest prefer{1]
[1 [n]1.19{]}[phr]Instruct me, {phr} [phr]for Thou know'st; {phr} [phr]Thou from the first{1]
[1 [n]1.20{]}Wast present, {phr} [phr]and, {phr} [phr]with mighty wings outspread, {phr}{1]
[1 [n]1.21{]}[phr]Dove-like sat'st brooding on the vast Abyss, {phr}{1]
```

Much like XML except

1. Relax well-formedness constraint on element typing (naming) in tags (“last in, first out”) while retaining tag types (start, end, empty)
2. Refactor and generalize attributes into annotations
3. Do without DTD (or top-down tree-based) validation/parsing for now
4. Tweak delimiters

Recognizably not XML any more

Does not conflict with XML if either is embedded in the other

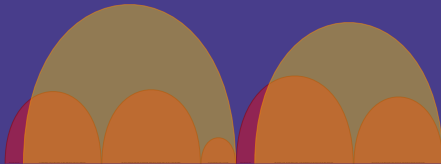
Note: this is nothing against XML!

(On the contrary it will show how useful and powerful it is)

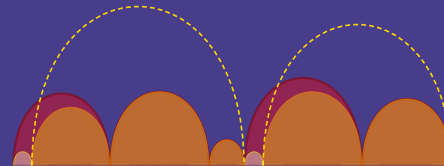
If XML didn't exist we'd have to invent it

In particular, XML turns out to be
a very capable platform
for implementing LMNL

LMNL supports overlapping ranges



In XML we fake it



This is the guy who once told Mike Kay that XSLT is fun


Neither is it a pitch (except in a general sense)

Building a **whole new stack** is a lot of effort!

Plus, there are many **other ways to deal** with overlap and annotation
LMNL (as a useful thing, not just a rumor) **would come with its own problems**

(If history is any guide)

Please continue to **use XML, as I will** (or your own favorite)

A hand is holding a red wax seal. The seal is circular and has embossed text that reads "Representation of Text". The background shows a bookshelf with several books. One book spine is visible with the text "EGYPTIAN LITERATURE". Another book spine has "SHELL POETRY". A third book spine has "WRITINGS". A fourth book spine has "RMIN". A fifth book spine has "PROSE OF VINCE".

Think of LMNL as an illustration, object lesson, or toy
Or as (a) material, a medium for making ...

LMNL Applications

Query Analysis Transformation

In theory

Because ranges in LMNL documents have no hierarchy, any (all) hierarchies are implicit

Any hierarchy may be interpolated when needed (and expressed in XML)

Arbitrary overlap (ranges overlapping ranges of the same type or family) is fine

Heterogeneous document description and annotation is possible

So is dynamic (auto-) tagging (prior markup does not interfere)

Visualizations and heuristics can expose structures or lack thereof

In practice

This is essentially a solo research project ...

(so progress is fitful)

... serving mainly as an opportunity for reflection ...

... *and yet* ...

Finally, pudding

Until 2011, my explorations all used XML (milestones) to encode LMNL

This is cumbersome and difficult, even while you learn a lot about processing XML

In the meantime I learned not to worry about “overlap” so much (as XML tools have improved)

Eventually (after others had done the same) I figured out how I could process the syntax

(LMNL syntax parsers or processors had already been implemented

by Gavin Thos. Nichol, Jeni Tennison, and Matt Palmer at least)

By thinking of parsing as not a series of events

But a sequence (pipeline) of mappings from syntax to (emerging) object model

Effectively, using a series of XSLT “refinements” to “compile” LMNL syntax

Building an XML-based (standoff) representation of a LMNL document

(Suitable for further processing)

Luminescent

XML into LMNL

LMNL to XML

XSLT, under (your choice)

Download at github.com/wendellpiez/Luminescent

Viewable at cocoon.lis.illinois.edu:8080/lis590dpl/wapiez/Luminescent/lmnl

Cocoon

LMNL to HTML/SVG

XProc (Calabash)

Heuristics / Analytics

XQuery (BaseX)

Lessons of LMNL

XML excels when you need a high degree of consistency,
especially over large amounts of information
with the disadvantage of having to define everything up front

In comparison to XML, LMNL is flexible and permissive
and correspondingly touchy in production (lacking validation guardrails)

But well-suited for small- and medium-size experiments including one-offs

Markup first, then model

Potentially at scale too (given better tools and interfaces)

Plus, LMNL offers a variety of interesting possibilities

still unexplored

Things difficult in XML are easy in LMNL

Things easy in XML are difficult in LMNL

A solution: use both together?

Complexity of the problem domain does not disappear, but it moves

Markup can be simple and spare, or more complex, ornate and promiscuous

— For complex tagging we need a structured editing UI

Some applications become truly easy

The cost has been in the time to build tooling

This presents the developer with a dilemma:

Build generic utilities, or focus on particular demonstrations?

References and resources

Layered Markup and Annotation Language

(Historical) specs at www.lmnl-markup.org

With grateful acknowledgement to Jeni Tennison, Gavin Thomas Nichol, John Cowan, Steve DeRose, Alex Czymiel, Paul Caton, Matt Palmer, Gregor Middell (and owed to many others) for all their efforts

Along with the collegial support of Sperberg-McQueen, Huitfeldt, Witt, Durusau, Durand, and many more

Luminescent

XSLT pipeline for parsing LMNL syntax into xLMNL, plus processing (in XSLT)

Capable of filtering, querying, XML extraction, display (HTML and SVG)

Requires XSLT 2.0; runs inside XProc (Calabash), Apache Cocoon, or BaseX (so far)

Open source at github.com/wendellpiez/Luminescent

Running on Cocoon at cocoon.lis.illinois.edu:8080/lis590dpl/wapiez/Luminescent/lmnl