

Tree Notation: an antifragile document notation

Breck Yunits

Abstract—A new minimal notation is presented for encoding tree data structures. Tree Notation may be useful as a base notation for domain specific languages. A new conjecture on data structures is introduced and predictions are made.

I. INTRODUCTION

Many widely used applications read and write documents in a domain specific language (DSL) based on XML,¹ JSON² or Racket.³ This paper and accompanying source code (github.com/breck7/treenotation) introduce a new whitespace-based notation that can serve as an alternative base notation for DSLs. These DSLs, called ETNs ("Extends Tree Notation"), are easy to create and can be very powerful. This paper describes Tree Notation (TN), three benefits and three downsides. It then introduces a new conjecture arising from TN and makes some falsifiable predictions.

II. TREE NOTATION

TN encodes one data structure, a **TreeNode**, with two members: an array of strings called **words** and an optional array of child TreeNodes called **children**.

TN defines three special characters, which in the canonical notation are node delimiter ("`\n`"), node edge ("`\"`"), and word delimiter ("`\"`"). The latter two are single spaces, not tabs.

A comparison quickly illustrates nearly the entirety of the notation:

JSON:

```
{
  "title" : "About Ada",
  "stats": {
    "pageViews": 42
  }
}
```

Tree Notation:

```
title About Ada
stats
  pageviews 42
```

III. PRACTICAL ADVANTAGES

A. Easy composition

When a document is composed of blocks written in multiple languages, those blocks may require verbose encoding to accommodate the underlying base notation.

In the multi-lingual example below, a JSON-backed IPython notebook encodes Python to JSON. The resulting document is more complex:

Breck Yunits is a software engineer in Seattle (breck7@gmail.com).

```
{
  "source": [
    "import matplotlib.pyplot as plt\n",
    "import numpy as np\n",
    "print(\"ok\")"
  ]
}
```

With TN, the Python block is indented and requires no additional transformation:

```
source
  import matplotlib.pyplot as plt
  import numpy as np
  print("ok")
```

In this example, the ETN would define a "source" node which interprets its child nodes as a single string.

B. Semantic diffs

JSON, XML, and Racket serializers can encode the same object to different documents by varying whitespace. Although ignoring whitespace can be a desirable feature in a language, programs can serialize identical objects differently, which often causes large diffs—and sometimes merge conflicts—for small or non-existent semantic changes.

In TN, editors have only one way to serialize an object. Diffs contain only semantic meaning.

C. Zero parse errors

Parse errors do not exist in TN. Every document is a valid TN document. Errors can only occur at the ETN level (e.g., a mistyped property name). Typos made at the location of a TN syntax character affect only the related node(s).

With other base notations, to get from a blank document to a certain valid document in keystroke increments requires stops at invalid documents. With TN all intermediate steps are valid.

A user can edit the nodes of a document at runtime with no risk of breaking the parsing of the entire document. If a node contains an ETN error, an editor can substitute an error handling node type for quick and/or automated correction.

A developer working on an editor that allows a user to edit the document source does not have to worry about handling both errors at the DSL level and errors at the base notation level. TN eliminates the latter class of errors.

IV. DRAWBACKS

A. Lack of Tooling and Support

TN is new, and library and application support, compared to other popular base notations, rounds to zero.

B. Lack of Primitive Data Types

Some notations, such as JSON, specify notations for common primitive types like booleans and numbers. Encoded documents can then be parsed directly to the matching in-memory structures. TN is relatively minimal and delegates the encoding and decoding of additional types to ETNs. Thus, parsing a TN document into the desired in-memory structures often requires the specification of an ETN.

In practice, however, rarely are the primitive data structures in a base level notation like JSON enough to fully parse a structure, and usually an additional parse step is used.⁴

C. Aesthetic Differences

Some developers dislike whitespace notations. In addition, without an ETN, a complex node in TN may be verbose, extending over multiple lines, with one node per line, whereas other base notations may benefit from a denser display of information out of the gate, with multiple nodes per line.

V. THE TREE CONJECTURE (TTC)

Working with TN leads to some observations about the general nature of data. To the LISP programmer, all data structures are recursive lists (or perhaps "recursive cons"). With the introduction of TN, this general notion can be stated in a specific way: **All data structures are trees.**

A. Steps to Disprove TTC

To disprove TTC, one needs to simply find a single instance of a structure that can't be serialized to TN via an ETN.

ETNs are subsets of TN. There exists an infinite number of ETNs. For example, a Hash Table ETN might serialize an instance to:

```
c 106
b 226
```

To the ETN Interpreter, this serialization represents a hash table, while at the same time to the TN interpreter it is a tree. Therefore, hash tables are a type of tree.

In a scan of over 300 programming languages and notations from over 50 years, no structure has been found that cannot be directly encoded by an ETN. Therefore, it is conjectured that all data structures are trees.

B. Predictions from the Tree Conjecture

TTC leads to many predictions. Two of them are listed below.

Prediction 1: Useful new ETNs may be found for every programming language.

Below is an example instance encoded in JTree, a JSON ETN.

```
o
s: title About Ada
o: stats
n: published 6182012
```

This ETN is one of the many languages in the vast TN universe. Some simple ETNs have been around for decades.⁵ ETNs are not limited to declarative data notations. There exists ETNs for C, Clojure, Ruby, et cetera, awaiting discovery. As soon as a new ETN is designed or found, any general ETN editor can then edit documents in that new ETN.

Prediction 2: The simplest 2D text encoding for neural networks will be an ETN.

By TTC, all trained neural networks are trees and therefore have ETNs. High level ETNs will be found to reduce these large trees into understandable notations.

VI. DISCOVERY PROCESS AND FUTURE RESEARCH

TN was discovered by an iterative process, where the goal was to find a useful language without violating certain natural Cartesian geometrical constraints. Perhaps an analogous process could be used for finding notations in higher dimensions.

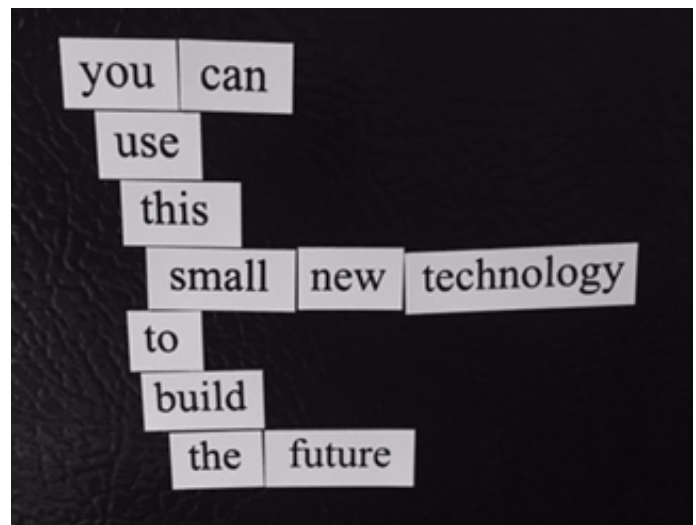


Fig. 1. Rearranging these fridge magnets is equivalent to editing a TN document. The fridge magnet set that includes parentheses is a poor seller.

VII. CONCLUSION

TN is a useful base level notation for DSLs. TN supports clean multi-lingual composition, aligns well with version control paradigms, and has an unbreakable base grammar that allows for robust runtime editing.

With more tooling support and the discovery of more ETNs, TN will improve developer productivity and enable faster and easier creation and maintenance of software.

TTC also makes certain predictions about future discoveries, and provides a basis on which to make further predictions.

REFERENCES

- ¹ Bray, Tim, et al. "Extensible markup language (XML)." World Wide Web Consortium Recommendation REC-xml-19980210. <http://www.w3.org/TR/1998/REC-xml-19980210> 16 (1998): 16.
- ² Crockford, Douglas. "The application/json media type for javascript object notation (json)." (2006).
- ³ Tobin-Hochstadt, Sam, et al. "Languages as libraries." ACM SIGPLAN Notices. Vol. 46. No. 6. ACM, 2011.
- ⁴ Ooms, Jeroen. "The jsonlite package: A practical and consistent mapping between json data and r objects." arXiv preprint arXiv:1403.2805 (2014).
- ⁵ Roughan, Matthew, and Jonathan Tuke. "Unravelling graph-exchange file formats." arXiv preprint arXiv:1503.02781 (2015).