

1 Making Open Transportation Data Useful and
2 Accessible: Recommendations for Good Practices in
3 Open Data Standards Management.

4
5 Elizabeth Sall*
6 UrbanLabs LLC
7 4405 4th Ave NE Seattle WA 98105
8 919 302 0265
9 esasall@gmail.com

10
11 Lisa Zorn
12 Metropolitan Transportation Commission
13 375 Beale Street, Suite 800 San Francisco CA 98105
14 415 778 6644
15 lzorn@mtc.ca.gov

16
17 Drew Cooper
18 San Francisco County Transportation Authority
19 1455 Market Street 22nd FL San Francisco CA 94103
20 415 522 4800
21 drew.cooper@sfcta.org

22
23 Bhargava Sana
24 San Francisco County Transportation Authority
25 1455 Market Street 22nd FL San Francisco CA 94103
26 415 522 4800
27 bhargava.sana@sfcta.org

28
29 Stefan Coe
30 Puget Sound Regional Council
31 1011 Western Ave #500 Seattle WA 98104
32 206 464 7090
33 scoe@psrc.org

34
35 Word Count: 6,773 + 250 = 7,023

36
37 Submitted to the 2017 Transportation Research Board Annual Meeting in response to the Standing
38 Committee on Library and Information Science for Transportation's Call For Papers: [*Innovative*](#)
39 [*Transportation Information & Data Tools and Practices*](#)

40
41 Submitted August 1, 2016

1 ABSTRACT

2 Effective data standards are critical to ensuring that open data is useful and accessible to its intended
3 audience. However, historical standards making organizations are using processes that aren't agile
4 enough, accessible, or relevant to the variety and quantity of data that is being introduced from the
5 multitude of recent open government data policies and pushes. This paper examines a variety of methods
6 for creating and managing open data standards to ascertain best practices. After evaluating the
7 applicability of these best practices to a real life standards-developing need, the project team concluded
8 that there were significant leadership needs in the data standards for transportation analysis arena.

9 INTRODUCTION

10 While open data policies and advocacy has resulted in an explosion in the number of available datasets
11 across a variety of sectors, some groups in the open data movement have moved beyond an initial push to
12 get as much data online as quickly as possible to increase *access*, towards measuring success by how open
13 data is *used* to support answering questions to solve a particular problem (1). Going one step further,
14 researchers have been investigating not just if open data is used, but *by whom* (2) and have suggested the
15 creation of a new digital-divide between those who have the skills and infrastructure to make use of the
16 data for their benefit and those who don't. These researchers have suggested that in addition to free and
17 open data, that quality documentation, technology access, and technical assistance must also be free and
18 open in order to empower those not already empowered. These issues can be broken into supply- and
19 demand-side issues, with this paper focusing on supplying the data in such a way that it is truly *open* to
20 the most potential users.

21 Most open data policies and manifestos published by local and federal governments as well as
22 advocacy groups list factors in making the open data supply *accessible* (3). Some of these are
23 straightforward and well defined such as requirements for machine readable formats and permissive use
24 licenses. Less well defined or mature are mandates for meta-data and the use of *standards maintained by*
25 *a standards organization* (4-5). The Office of Management and Budget's Open Data Policy (created in
26 response to President Obama's Executive Order (6)) points to a "common-core metadata schema" that is
27 further defined on its website. A *data standard* is a set of rules on the format and meaning of data to
28 facilitate its sharing and exchange. The Open Knowledge Foundation, an organization devoted to
29 promoting the free exchange of information defines an *open data standard* as one that is free from
30 licensing restrictions and developed to be vendor-neutral (7). After giving a little bit of background on
31 who makes standards and how we might be able to evaluate them, this paper analyzes several case studies
32 of data standard development and management, recommends a set of best practices from them, and
33 evaluates how implementable the best practices are with a new data standard creation project.

34 Background: Who Makes Data Standards?

35 So who makes all of these standards that these recommendations discuss, and how do we make sure they
36 are developed and managed in a way that will make sure they remain useful to the appropriate audiences?
37 Starting from the top-down, there is the International Standards Organization (ISO), the American
38 National Standards Institute (ANSI), a variety of other "accredited" standards development organizations
39 (SDOs), "non-accredited" official standards development organizations, companies, and ad-hoc, or loose
40 standards. The Federal Government has had limited, but strategic involvement in standards-making per
41 *OMB Circular A-119: Federal Participation in the Development and Use of Voluntary Consensus*
42 *Standards and in Conformity Assessment Activities* (8) and clarified in the 2012 OMB Memo M-12-08
43 *Principles for Federal Engagement in Standards Activities to Address National Priorities* (9) which

1 reaffirms that the development of standards be left to the private sectors **but for** instances of national
2 priority identified in law where government may take on the role of the catalyst to help progress a market-
3 based solution.

4 The most recognized body for standards is ISO, a non-governmental body that develops standards
5 within over 250 different technical committees. Currently, the only standards related to the transportation
6 sector are related to vehicle communication, intelligent transportation systems, and shipping. ISO has
7 four key principles in standards development: responsive to the market need, based on global expert
8 opinion, participation from multiple stakeholders, and based on a consensus (10).

9 ANSI is a member organization of ISO, but has the specific mission of progressing the US
10 economy by promoting and progressing standards. ANSI accredits the processes of over 200 standards
11 developing organizations (SDOs) that have created over 10,000 standards. ANSI requires that standards
12 be developed by consensus of those affected by them, in a public manner, and with due process.

13 (11)

14 Almost all of the accredited SDOs (12) are industry-specific groups or professional organizations
15 such as the American Society of Civil Engineers. All of the two-dozen ANSI-accredited SDOs that the
16 team reviewed as a part of this project conducted their activities via static web-postings of PDFs or in a
17 non-public webspace, accessible by request. While technically fulfilling the definition of “open”, they
18 are not conducive strategies to encourage dynamic dialogue or friendly challenges that would be the
19 hallmark of an agile, open standard.

20 The federal government does sponsor a number of standards-creation or endorsement activities in
21 order to advance items of strategic importance. This includes the Federal Geographic Data Committee
22 (13) which both makes and endorses standards.

23 Finally, there is no shortage of “open” data standards, successful and not, developed at companies
24 (e.g. GTFS at Google) or by individual projects or unaccredited organizations (e.g. Open Street Map).

25 **METHODOLOGY**

26 The approach to developing recommendations regarding open data standards is based on evaluating a set
27 of case studies both in and adjacent to the transportation industry to glean a set of best-practices. This
28 paper evaluates how implementable these best practices are within the context of a project from which we
29 derive a series of recommendations and conclusions.

30 In order to identify characteristics and practices of successful data standards, the team first identified a
31 handful of measures that would separate out successful from unsuccessful data standards and then
32 evaluated a handful of examples from the transportation industry and beyond.

33 Measuring the effectiveness of a data standard could be done in a myriad of ways. One way, is to
34 measure the *market share*, or the degree to which the standard **is** used when it **could** be used. Another, is
35 how *compatible* the standard is with other popular data schema and standards that operate in the same
36 ecosystem. *Rigor*, whether a standard is clear and well defined, is important because otherwise isn't much
37 use and could result in real damage. We have infamously sent rockets on the wrong trajectories because
38 of non-rigorous data standards. The next set of methods for measuring a useful data standard are likely to
39 be more qualitative in nature. *User satisfaction*, or the degree of enthusiasm with which people use the
40 standard can be used to distinguish a much despised, but entrenched standard from a sparsely used, but
41 beloved one. *Approachability*, or the level of information and technical pre-requisites, can help address
42 concerns about increasing the digital divide.

1 Depending on what a standards is used for and its *audience*, some of these measures will be more
2 important than the other. For the purposes of this paper, we define the following audience types:

- 3 • Technical: completely backend and never intended for public consumption and use. Market
4 share, compatibility, and rigor are the most important measures here. Example: Vehicle data,
5 USB Specs.
- 6 • Technical Consumer: approachable to people with no specific industry training, but general
7 technical abilities. Approachability becomes more and rigor slightly less important measures.
8 User satisfaction becomes more important because people could easily defect to another standard
9 of their choosing which would cause market share and interoperability to decline. Example:
10 GTFS, HTML.
- 11 • Public Consumer: hobbyists with no particular training can easily understand and use. Measures
12 other than approachability are much less important as the datasets are usually simple enough and
13 there is little likelihood that they need to interface rigorously with other data. Example: Street
14 Tree List

15 Because there is less rigor required for public consumer-oriented data standards, the discussion and case
16 studies will focus on technical and technical consumer-oriented data standards.

17 **CASE STUDIES**

18 The following cases span industries, audiences, and levels of success and touch upon not just the
19 development of the data standard, but its management. Some of the cases pertain to specific data
20 standards while others identify how an organization responsible for multiple data standards operates.

21 *Universal Traffic Data Format (UTDF) - Not Open Enough*

22 Trafficware created the UTDF to import and export data to and from its Synchro traffic engineering
23 software (14). Synchro and the format are widely used by traffic engineers and the format has seen
24 support from a variety of other open source projects such as DTALite (15). However widely used, it is
25 not publically managed or even documented and is thus not an “open” format in the common use of the
26 term.
27

28 *TransXML - Lack of Persistent Management*

29 The TransXML standard was developed as cooperative research program project NCHRP 20-64 and
30 published as NCHRP Report 576 (16). While the specifications developed as a part of this project are
31 technically open, there does not appear to be widespread use of any of the formats put forward. There are
32 three primary reasons that this specification failed to be popular. First, it was developed using a typical
33 contracting relationship rather than a truly collaborative process involving the actual users of the standard.
34 Second, there is not any evidence of continued support to embrace the standard: the last update to the
35 TransXML project website was in 2006, shortly before the official NCHRP project concluded. Finally,
36 the standard itself wasn't very accessible. It was buried in a large, less searchable PDF report with no
37 code validator, example data, or copy-pasteable text.
38

39 *General Transit Feed Specification (GTFS) – Outgrowing its Original Purpose*

40 GTFS (17) was created by and is managed by Google Inc., which invested a fair amount of resources in
41 making it free and very easy to convert existing proprietary transit service formats to GTFS, to make it
42 easy for Google to use the GTFS data in their maps application. Google used a very accessible file format
43

1 that could be opened up and read by just about everybody with a computer (CSV), and a vocabulary that
2 made it extremely human-readable at the expense of computationally efficient (e.g., “bus” rather than
3 some integer code that means bus, and “Mission St/24th St” rather than some code). Their format is well
4 documented on their website and there is a Feed Validator that allows users to view the files in their maps
5 application to check for errors.

6 Despite GTFS being very ubiquitous, its user satisfaction is not as high as it could be. One reason
7 could be that while declaring itself “not set in stone”, there is very little that Google does to actively
8 manage the format to respond to new and growing needs.

9 In a sense, GTFS is a victim of its own success. Its ubiquity and availability has opened it up to
10 new and growing audiences who need more and more from the standard. And Google, this being outside
11 of their main line of business, has not adequately responded to the demand. Changes to the specification
12 are requested through an email listserv without an obvious release structure or version control system
13 outside of an on-going change-log. While GTFS has declared itself to always remain backwards-
14 compatible (a double-edged sword), version control is perhaps paramount in a good data standard. Oddly,
15 Google does have its GTFS Feed Validator under version control, but the [Feed Validator](#) does not use
16 GTFS version (because there is none) as an input.

17 *Open Street Map (OSM) - Too much help, not enough central power*

18 OSM’s XML-based format uses a wikimedia-flavored wiki for documentation, which is both extensive,
19 easy for editors to update, and easy for editors and users alike to view the history/changelog of the page
20 (18). The data vocabulary is well documented with examples. The wiki notes the various pros and cons
21 of using an XML format: while it is fairly human and very machine readable, it is not space-efficient. As
22 a space-saving, but illegible alternative, OSM also provides a Protocolbuffer Binary Format (PBF). There
23 are two main deficiencies with the OSM data standard itself: (1) its lack of interoperability with existing
24 standards or data vocabularies; (2) technocratic nature of the membership has led to, despite their best
25 intentions, an “ungrokable” data set by people with a normal level of technical skill.

26 The OSM data format is managed by the OSM Foundation, a bare-bones, mainly volunteer effort.
27 Despite their best efforts, the Foundation’s underwhelming and confusing primary support is often
28 trumped by superfluous and conflicting secondary documentation, often provided by a for profit
29 company. The OSM Foundation allows for a well-documented, formal public [proposal process](#) for
30 adding tags and relationships, but it is clunky to navigate and the open member voting allows for people
31 to reject things with unhelpful comments such as “[history is just plain wrong](#)”. It would be more helpful
32 to have a smaller, dedicated group of people who are vested in the standard to be the primary stewards of
33 any changes. Finally, while these changes seem to be much more dynamic and responsive to user needs
34 than GTFS, there is no version control of the standard to speak of.

35 *Traffic Data Exchange Format - Fragile but exciting*

36 A consortium of not-for-profit entities (The World Bank) and companies who had a history of
37 collaborating together on open transit data using GTFS sought to create an open data standard that would
38 allow them to display traffic data in map tiles (19). While well documented with examples and managed
39 and available on GitHub, the specification itself is too specific and hasn’t yet considered a variety of use
40 cases. It is also fragile as it is overly dependent on specific revision numbers in OSM data, which
41 exposes a flaw not just in this format, but in OSM’s as well. While the collaborative effort to create a
42
43

1 data standard for transportation performance measures is promising, it needs some more thought and
2 investment in order to take root.

3 *OpenTrails*

4 OpenTrails is a project incubated by Code for America to define a geojson-based data structure for
5 recreational facilities (20). It is new, but has a fair amount of traction and a good consortium of partners
6 which gives a good indication that it will soon have substantial market share over the existing spread of
7 various GIS formats residing on people's desktops and an existing federal trail standard that is too
8 narrowly defined and buried deep in government-controlled PDF documents. OpenTrails is well
9 documented via its website in an approachable format and allows the public to comment and make
10 suggestions to its specification via a google document. Unfortunately, while these comments are likely
11 read and incorporated to some degree, all decisions about the specification are made in a temporally-
12 sorted listserv, making it difficult to look up a change-request for a particular field. There is also no way
13 to compare versions of the standard dynamically or even via a changelog.

14
15 One useful tactic that the OpenTrails team used was to have a defined "Request for Comment"
16 period between version 1.0 and 1.1 and then a freeze on development for a period of time after version 1.1
17 was released to allow users update their data. OpenTrails, similar to GTFS, also publishes a data
18 validator and example data. However, they are confusingly maintained in the same Git repository as the
19 main OpenTrails website.

20 21 *W3C - Necessary for experts; ungrokable for newcomers.*

22 The World Wide Web Consortium (W3C), the organization that standardizes the bulk of the soft internet,
23 develops standards within working groups. While most of the W3C work takes place within a mediaWiki
24 instance, each working group undertakes the task of discussing and iteratively developing the standard
25 differently, but there has been a trend recently to shift the development to a public GitHub repository. As
26 an example, the Spatial Data on the Web Working Group (21) maintains their normal MediaWiki site, but
27 also has a GitHub repository (22) where their standards are actually worked out and documented. W3C
28 maintains an "official site" for each standard along with links to the previous or new versions if
29 applicable. They number the release of each standard (e.g., HTML 5.0, 5.1, etc) but use dates for small
30 changes made between such as HTML 5.1 from September 29th 2015. Perhaps most interesting is their
31 use of a strict framework, the Resource Description Framework (RDF)(23), for defining relationships and
32 ontology and a language, Turtle, that is capable of defining the RDF in plain-text. The biggest drawbacks
33 from the W3C process are (A) the variety in ways in which the standards are developed makes it difficult
34 to understand how to participate in the process in a meaningful way, and (B) the process and the standards
35 themselves are designed for experts. They are complicated to comprehend.

36 In addition to its work on standards, W3C has developed a draft *Data on the Web Best Practices*
37 based on a number of case studies researched by the Data on the Web Working Group. Best practices
38 pertaining to data standards include recommendations to use multiple, open, machine-readable,
39 standardized data formats to promote interoperability, and to reuse existing data vocabularies when
40 possible.

41 42 *Internet Engineering Steering Group (IESG) - Useful and Stylish Documentation for Expert Technicians*

43 IESG manages the internet's engineering standards process, which dates back to the days of DARPA net.
44 The process for developing internet standards is documented by its own standard: RFC2026(24).

1 Standards are developed through the “Request for Comment” (better known as RFC) process. All RFCs
2 are maintained on www.rfc-editor.org and are assigned a serial number. Not all RFCs are standards, but
3 those that are go through an evolution from *proposed* -> *draft* -> *standard*. Standards each have serial
4 numbers (e.g. [STD 77](#)) and new RFCs can update or make old standards obsolete. So the STANDARD
5 has a number and the RFC number is changed to reflect changes to that standard.

6 RFC-editor maintains a rolling list of status-changes. The problem is that if you are using a
7 canonical reference to the RFC, you might not ever know that anything has changed or moved forward.

8 There is a very popular RFC “style guide” (25), which is used for defining many standards across
9 many disciplines, not just IESG. IESG also allows another type of RFC called “best current practice”
10 (BCP), which are not standards, but define the coalescence around a certain way of doing things in
11 practice.

12 *Frictionless Data*

14 Frictionless Data is a project of the Open Knowledge Foundation, a not-for-profit dedicated to making
15 knowledge accessible and making sure people are able to use it (26). One of their key contributions is the
16 definition of a Data Package including a meta-data standard. One of the most interesting things about
17 these standards is that they are developed and published on the web using a GitHub repository that
18 seamlessly creates a Jekyll-based webpage based on an RFC-style specification that lives in a version-
19 controlled markdown file. Because the entire operation is in GitHub, you can easily access a change-log
20 and do key differences. Moreover, it is very human-readable and easy to interact with using GitHub’s
21 built-in markdown renderer, issue tracking system, and commenting system. It is also easy to use when
22 downloaded *from* GitHub, since it just uses ASCII text.

23 This approach is nice for several reasons. First, it forces any change to be noticed. Because you
24 are using version control software, you can’t ‘cheat’ and slip in something else to a previous version. It is
25 ASCII-based, which allows for easy viewing of differences across versions of the standard. It is human-
26 readable and allows for more formatting than the pure ASCII view of RFCs. It has a seamless translation
27 between the version-controlled ASCII and the front-end webpage which reduces possible error or
28 misalignment. Finally, it leverages GitHub’s toolset which will allow you to subscribe to a repository
29 which will let you keep abreast of any changes.

30 **DISCUSSION**

31 These case studies have a broad range of success because of or in spite of their origin story and continued
32 management strategy. Despite the broad number of “best practices” in open government data
33 management that point to using defined, open, and interoperable standards, there is little detail aside from
34 broad ideals that pertain to the development and management of these standards. Additionally, because
35 the vast amount of data collected at various levels of government does **not** likely already have a relevant
36 standard to point to, there needs to be a large increase in the number of standards successfully developed
37 and maintained in order to meet this need.

38 The primary takeaways for technical standards from a review of these case studies are:

- 39 1. Developed and managed by a trusted source with permanence, customer focus, and with sufficient
40 user involvement
- 41 2. Leverage existing data vocabularies
- 42 3. Right-size the standard (and its management) to the audience
- 43 4. Evolve the standard at right pace and using rigorous methods

- 1 5. Limit unnecessary tools and libraries
- 2 6. Diligent documentation of standard and process
- 3 7. Balance flexibility while limiting vocabulary dispersion
- 4 8. Structure and tools that limit and catch errors
- 5 9. Promote your standard to make sure industry knows it is there

6 The next section evaluates how implementable these recommendations are by trying to implement them
7 when designing a new data standard based on a project need.

8 **IMPLEMENTATION EVALUATION**

9 When three public agencies, the Metropolitan Transportation Commission (MTC), Puget Sound Regional
10 Council (PSRC), and San Francisco County Transportation Authority (SFCTA) decided to implement a
11 user-ready version of some software developed at a University, one of the first tasks was to develop a set
12 of interoperable data standards that could ensure that all three agencies (plus others) would be able to use
13 the software. The software, Fast-Trips (27), does person-based dynamic transit passenger assignment and
14 requires passenger demand and transit networks as inputs.

15 This section evaluates how this interagency team considered the case study findings to develop
16 both the GTFS-PLUS transit network (28) and the dyno-demand demand (29) data standards and
17 evaluates the degree to which the case studies findings are implementable. These standards were
18 primarily developed for a purely *technical* audience, but envisioned to be useful to a broader, *technical*
19 *consumer* audience thus accessibility was highly valued.

20 **Developed and Managed by a Trusted Source with Permanence, Customer Focus, and with** 21 **Sufficient User Involvement**

22 The first problem the Agency Team faced was that there was a need for a standard “as soon as possible,”
23 but there was no relevant, existing SDO with the technical domain expertise within which to officially
24 conduct the business of standard development. The Transportation Research Board (TRB) has relevant
25 domain expertise among its volunteers and has previously directed standards-development activity, but
26 the timeline for this was years, not months. The American Society of Civil Engineers is the closest SDO
27 to this domain, but their process was not conducive to the timeline or technical needs of the project. On
28 the transit networks side, the team considered trying to make changes to the official GTFS specification,
29 but that process as previously noted is ad-hoc and many of the variables that the project needed did not fit
30 within the requirements that Google has put forth for new variable names.

31 As a workaround to these issues, the team developed the data standards publically under the
32 auspices of the “Open Source Planning Data Standards” GitHub handle and worked to fully document the
33 standard beyond the needs of the project.

34

35 *Findings:*

- 36 ● Need more SDOs (accredited or not) to cover areas of expertise where open data is expected to be
37 released, but not too many such that it isn’t clear who is responsible.
- 38 ● Need a directory of SDOs by domain and a standardized process for “involving appropriate
39 users” so that standards aren’t re-developed.
- 40 ● Standards creation activities need to be responsive to the timeline of project needs: weeks, not
41 years.

42 *Workaround:*

- If you need to develop the standard within the context of a project without a home organization, separate the standard as much as possible from the project and prepare to be flexible and evolve the standard to other users' inputs.

Leverage Existing Data Vocabularies

Given the general industry coalescence around GTFS, the team decided to start with its standard and data vocabulary in order to maximize the overlap in required data and leverage existing tools. While GTFS contains a significant amount of data, it doesn't have all the detail or structural relationships required for our project. However, the changes required would both take months to years to make within the existing GTFS management process as well as break the rules for GTFS to be backwards compatible.

At the outset, the team decided to make everything that was mandatory in GTFS mandatory in the new standard even if it wasn't needed to run Fast-Trips. Starting with GTFS meant an implicit adoption of a CSV-based format with a header line containing case-sensitive variable names. In addition to file type, in order to utilize existing GTFS validator and visualization tools, the existing GTFS files needed to stay intact and not add unexpected fields. Therefore, the team adopted a transit network standard that had additional files with supporting information.

Generally speaking, for every GTFS file, there is an additional file with the additional fields as well as a identifier key to link it back to the GTFS information. For example, `stops.txt` would be a standard GTFS file, but and `stops_ft.txt` contains a variable `stop_id` to link back to `stops.txt` as well as the additional information needed for Fast-Trips such as the presence of shelter, lighting, seating, etc.

Finding:

- Existing data standards that the industry has coalesced around may not be perfect.

Workaround:

- Don't throw the baby out with the bathwater. Reuse existing data vocabulary and software by "extending" the standard rather than creating a new one.

Right-Size the Standard (and Its Management) to the Audience

In the Fast-Trips case, the team was primarily developing this set of standards for a purely technical audience of professionals in the travel analysis industry with the technical public a secondary audience. That said, the in-depth and rigorous standards development processes found at W3C, ISO, and the Open Street Map Foundation would likely scare away all but the most technical of the technical consumers. Nobody is likely to read a several-hundred-page specification document or read through a document written in the Turtle language to evaluate if the right variable ontology has been used. Thus, the team tried to strike a balance between documentation rigor and approachability and limited the introduction of new tools, skills, and terminology that people in the travel analysis industry didn't already use.

Finding:

- Rigorous and overly detailed processes and standards documentation would likely deter all but the most technical and determined individuals in our professional niche.

Workaround:

- Use tools and terminology already familiar to the industry.

1 Evolve the Standard at Right Pace and Using Rigorous Methods

2 During the first few months after an initial draft, the standard was in constant flux as the team put it
3 through the paces of use in the real world. This flexibility in the rapid evolution was necessary in order to
4 get a standard that worked for the purposes of the project but would be completely untenable to a broader
5 audience who would be trying to develop tools that could also interact with the standard.

6 In order to mitigate this, the team adopted semantic versioning (<http://semver.org/>) where
7 anything before version 1.0.0 should be considered in pre-release and subject to change at any time. Once
8 version 1.0.0 is released, any changes that break backwards compatibility should be contained to annual
9 or sparser events. Implementing strict version releases also allows rigor in defining requirements from a
10 software perspective (e.g., this software requires data be defined using standard version 1.2.0 or higher).

11 In order to enforce rigorous version control and leverage their existing cooperative development
12 tools, the team decided to put the standard on GitHub as Markdown-flavored ASCII text. The git version
13 control software that GitHub is based on enforces that every version is stored and allows for easy
14 comparison of differences between versions.

15 The collaboration tools available on GitHub such as issue-tracking (see Figure 1) and resolution
16 allowed the team to raise issues with the standard, discuss fixes, and quickly implement resolutions in the
17 standard itself. GitHub offers substantial benefits over MediaWiki, which requires substantial initial
18 setup to achieve the structure already found in W3C and Open Street Map Foundation. GitHub is also
19 free for public projects whereas MediaWiki requires a server that would incur ongoing maintenance costs.

The screenshot shows a GitHub issue page for "Column name change in drive_access_points_ft.txt #13". The issue is marked as "Closed" and was opened by user "bhargavasana" on Feb 9. The issue title is "Column name change in drive_access_points_ft.txt #13". The issue is closed and was opened by user "bhargavasana" on Feb 9. The issue has one comment from "bhargavasana" suggesting a change from "lot_long" to "lot_Jon". A second comment from "e-lo" agrees and mentions queuing the change until another update. The issue has a "fix for clarity" label, an "opportunistic" milestone, and a commit that referenced it on Feb 17. The issue was merged and closed in #16 on Feb 17.

20
21 Figure 1 Example of Issue Management in GitHub

22 Findings:

- 23 • Use semantic versioning and be explicit about pre-releases using version number less than 1.0.0

- GitHub provides appropriate built-in tools for writing, managing, and collaborating on data standards without the initial setup time and continued cost of using MediaWiki.

Limit Unnecessary Tools and Libraries

When considering file formats, the Agency Team seriously considered HDF5 which would offer a highly compressible format for the large demand files and SQLITE which has a substantial number of built in query capabilities. However, it was telling that all three agencies didn't have any one of these tools installed in their standard computing environment and thus the team adopted a CSV-based standard that could be read and edited in a normal text editor if desired, or easily read in using standard libraries to almost every computer programming environment. The CSV-based format has the added benefit of being able to version control the data itself.

Diligent Documentation of Standard and Process

Each standard is documented using RFC-style text in a GitHub repository. In the repository base is a README.md file that has the following sections:

- Version
- Date last updated
- Date created
- Authors
- Changelog
- High level Known Issues
- Files that MAY or MUST be contained in the data to comply with the specification, with links to descriptions of those individual files.

Within a subfolder /files are separate ASCII-based markdown files describing each of the mandatory or optional files using RFC-style text describing the format type, and required and optional variables.

Balance Flexibility while Limiting Vocabulary Dispersion

The whole point of having standards is to come up with a set of expectations about the information that will be provided. However, even between the three participating agencies, there were often three different methods for approaching problems which required differences in the required data. An example of this is park and ride lots and whether they should be included as a singular lot, a list of lots, or not as an input at all with the travel demand. The team managed this particular issue by including the park and ride lot as an optional field which can contain a list of lots, including a list that just contains a single lot.

Another issue was coming up with a common vocabulary to talk about the same thing that was recommended, but not required. For example, mode names could be light-rail or light_rail. Our standard suggests (but does not require) that you use light_rail so that hopefully, eventually, everybody will use the same vocabulary, but it won't prevent you from using it if it doesn't work for your application.

A final strategy in accommodating a multitude of needs was to make only the very basic required data "required" as part of the standard and everything else optional. This is a double-edged sword because while it will make the standard accessible to many agencies who want to use it and don't have the data to support it (not every agency has information about every escalator or bus stop bench), it will create a very inconsistent level of data.

Findings:

- Small differences in workflows can sometimes require significant differences in data structures

- Flexibility in mandatory fields will result in wide variety of data being available and limit the scope of interoperability

Workaround:

- Specify different fields as mandatory for different applications. i.e. dyno-demand, platinum-level could mean that all mandatory and optional fields are filled and gives you access to the broadest number of uses.

Structure and Tools that Limit and Catch Errors

Limiting the number of times something is duplicated will reduce the opportunity for conflict and limit the number of places that data needs to be edited. The only place that data should be duplicated is when it serves as a key between files. Sometimes, as in the case of dyno-demand, this means separating out data into key component files: person, household, and trips so as to not necessitate duplicating person and household attributes across every row in the trips file. While this requires a bit of manipulation on the software side, it reduces file sizes and the opportunity for error.

Limiting the propagation of errors was another key reason to decide to allow plain-worded string types for variables like travel mode. Using number-string lookups for common words like modes, vehicle types, and trip purposes would save a small amount of space and input read time. However, if one opens a file and sees "1, 88", one is much less likely catch an error compared to reading "commute, pogo stick" which would immediately alerts one that there could be an error.

Finally, the data standard should include a validator that reads in the data and lists both fatal and potential errors. Fatal errors means that the data does not comply with the minimum standards. Potential errors are a list of things that you should consider before proceeding. These might include unreasonably fast or slow buses on the GTFS-PLUS side, or values-of-time outside normal parameters on the demand side. While the Agency Team did not write its own validators within the standard (instead using Google's GTFS validator), it is something that should be considered in the future.

Promote Your Standard to Make Sure Industry Knows It Is There

Lastly, no data standard usefully exists in isolation. In order to let people know what the team was working on, it was broadcast to the relevant user groups like the Travel Model Improvement Program Listserv as well as published on our project blog (<http://fast-trips.mtc.ca.gov>) and presented at the Transportation Research Board Innovations in Travel Modeling Conference and at informal conferences such as Transportation Camp. The team has also worked with other agencies interested in using the standard. That said, the standards need to have a home and team behind it that last longer than the length of this particular research project, and thus a SDO for the travel analysis industry must be developed in the short term with collective buy-in from the broad industry.

CONCLUSIONS

Based on the review of case studies and the Agency Team's experience in developing and managing their own data standard, the following conclusions have been reached.

Leadership Needed

Open data standards are an important and required component of making open government data useful and accessible, but are currently more of an afterthought in the process. Rather, the data standard should pre-date the data release. This currently isn't happening for a significant amount of data in the transportation sector because there is a lack of official coalescence around a standards development

1 organization that has standards development in its facility, mission, and funding. Rather, there are several
2 ad-hoc standards developed, each without significant market share and limited use cases and several
3 others managed by corporate interests. While the government has purposefully (and rightfully) decided to
4 take a backseat role in standards development, there is a role for a third party leader in modern data
5 standards creation and management that needs to be funded, empowered, and filled. It should be
6 emphasized that having strong leadership is not at odds with having a user-driven process. Rather, strong
7 leadership is required in order to facilitate a user-driven process that is accessible to the appropriate
8 audiences

9 **Tools for Open Standards Development Are Now Easy to Use and Accessible**

10 The Agency Team found GitHub to be an ideal tool to manage the development and documentation of a
11 data standard and much preferable to the start-up effort, cost, and continued maintenance of MediaWiki.
12 GitHub (and other tools) are free to use for public projects and obviate any excuses for conducting
13 matters behind closed doors/email groups/ or in indecipherable PDF document dumps.

14 **Using Best Practices Will Increase the Usability and Accessibility of a Data Standard**

15 This paper outlined a set of data standards best practices based on various case studies. The standard itself
16 should use existing data vocabularies to promote interoperability, limit unnecessary tools and libraries, be
17 thoroughly documented, balance flexibility and unneeded dispersion, and limit data duplication.

18 New standards should be developed and managed by a trusted source with permanence, customer
19 focus and with user involvement. The standards developer should be aware of evolving the standard too
20 fast or too slow and use rigorous version control and validation methods to limit errors. Finally, a data
21 standard needs to be promoted and supported in order for the right audiences to know to use it and
22 understand how.

23 **ACKNOWLEDGEMENTS**

24 The authors would like to thank the Federal Highway Program SHRP2 Implementation Assistance
25 Program which funded this work as well as the Fast-Trips Implementation Management Team and project
26 manager: David Ory, Billy Charlton, Joe Castiglione, Mark Simonson, and Diana Dorinson. Finally, the
27 authors would like to thank Mark Hickman and Alireza Khani, the original Fast-Trips designers.

28 **REFERENCES**

- 29 1. Susha, Iryna, Åke Grönlund, and Marijn Janssen. "Organizational Measures to Stimulate User
30 Engagement with Open Data." *Transforming Government: People, Process and Policy* 9, no. 2 (May
31 18, 2015): 181–206. doi:10.1108/tg-05-2014-0016.
- 32 2. Kaasenbrood, Maaïke, Anneke Zuiderwijk, Marijn Janssen, Martin de Jong, and Nitesh Bharosa.
33 "Exploring Factors Influencing the Adoption of Open Government Data by Private Organisations."
34 *International Journal of Public Administration in the Digital Age* 2, no 2 (2015): 92. Doi:
35 10.4018/ijpada.2015040105.
- 36 3. Burwell, Sylvia, Steven VanRoekel, Todd Park, and Dominic J. Mancini. "M-13-13 – Memorandum
37 for Heads of Executive Departments and Agencies: Open Data Policy – Managing Information as an
38 Asset." <https://project-open-data.cio.gov/policy-memo/>.
- 39 4. U.S. Senate. 114th Congress. S.B. 2852, *OPEN Government Data Act*. Washington. Introduced April
40 26, 2016. <https://www.congress.gov/bill/114th-congress/senate-bill/2852/all-info>.

- 1 5. U.S. House. 114th Congress. H.B. 5051, *OPEN Government Data Act*. Washington. Introduced April
2 26th, 2016. <https://www.congress.gov/bill/114th-congress/house-bill/5051/>
- 3 6. Executive Order of May 9, 2013, Making Open and Machine Readable the New Default for
4 Government Information. [https://www.whitehouse.gov/the-press-office/2013/05/09/executive-order-](https://www.whitehouse.gov/the-press-office/2013/05/09/executive-order-making-open-and-machine-readable-new-default-government-)
5 [making-open-and-machine-readable-new-default-government-](https://www.whitehouse.gov/the-press-office/2013/05/09/executive-order-making-open-and-machine-readable-new-default-government-)
- 6 7. “Open Data Handbook,” Open Knowledge Foundation, <http://opendatahandbook.org/guide/en/>.
- 7 8. Office of Management and Budget, “Draft Revision OMB Circular A-119: Federal Participation in
8 the Development and Use of Voluntary Consensus Standards and in Conformity Assessment
9 Activities.” [https://www.whitehouse.gov/sites/default/files/omb/inforeg/revised_circular_a-](https://www.whitehouse.gov/sites/default/files/omb/inforeg/revised_circular_a-119_as_of_1_22.pdf)
10 [119_as_of_1_22.pdf](https://www.whitehouse.gov/sites/default/files/omb/inforeg/revised_circular_a-119_as_of_1_22.pdf)
- 11 9. Chopra, Aneesh, Miriam Sapiro, and Cass R. Sunstein, “M-12-08 Memorandum for the Heads of
12 Executive Departments and Agencies: Principles for Federal Engagement in Standards Activities to
13 Address National Priorities.” January 17, 2012.
14 <https://www.whitehouse.gov/sites/default/files/omb/memoranda/2012/m-12-08.pdf>.
- 15 10. “International Standards Organization Website,”
16 http://www.iso.org/iso/home/standards_development.htm.
- 17 11. “American National Standards Institute,”
18 https://www.ansi.org/about_ansi/introduction/introduction.aspx.
- 19 12. “American National Standards Institute List of Accredited Standards Development Organizations,”
20 [https://share.ansi.org/Shared%20Documents/Standards%20Activities/American%20National%20Stan-](https://share.ansi.org/Shared%20Documents/Standards%20Activities/American%20National%20Standards/ANSI%20Accredited%20Standards%20Developers/JUNE16ASD.pdf)
21 [dards/ANSI%20Accredited%20Standards%20Developers/JUNE16ASD.pdf](https://share.ansi.org/Shared%20Documents/Standards%20Activities/American%20National%20Standards/ANSI%20Accredited%20Standards%20Developers/JUNE16ASD.pdf).
- 22 13. “Federal Geographic Data Committee Website,” <http://www.fgdc.gov/>.
- 23 14. “Trafficware Website,” <http://www.trafficware.com/transferring-data-using-utdf.html>.
- 24 15. “DTALite Release Notes,” https://github.com/xzhou99/dtalite_software_release.
- 25 16. Ziering, Eric, Frances Harrison, and Paul Scarponcini, “NCHRP Report 576; XML Schemas for
26 Exchange of Transportation Data.” Transportation Research Board, 2007:
27 http://onlinepubs.trb.org/onlinepubs/nchrp/nchrp_rpt_576.pdf.
- 28 17. “General Transit Feed Specification Website,” <https://developers.google.com/transit/gtfs/>.
- 29 18. “Open Street Map Data Standard,” http://wiki.openstreetmap.org/wiki/OSM_XML.
- 30 19. “Traffic Data Exchange Website,” <https://github.com/opentraffic/traffic-data-exchange-format>.
- 31 20. “Open Trail Specification,” <http://archive.codeforamerica.org/specifications/trails/spec.html>.
- 32 21. “W3C Spatial Data on the Web Working Group,” https://www.w3.org/2015/spatial/wiki/Main_Page.
- 33 22. “W3C Spatial Data on the Web Working Group Github Repository,” <https://github.com/w3c/sdw/>.
- 34 23. “W3C Resource Description Framework,” <https://www.w3.org/TR/rdf11-concepts/>.
- 35 24. “IETF Request for Change 2026,” <https://www.ietf.org/rfc/rfc2026.txt>.
- 36 25. “IETF Request for Change 7322,” <https://www.rfc-editor.org/rfc/rfc7322.txt>.
- 37 26. “Frictionless Data Website,” <http://frictionlessdata.io/>.
- 38 27. “Fast-Trips Software,” <https://github.com/MetropolitanTransportationCommission/fast-trips>.
- 39 28. “GTFS-PLUS Transit Data Standard,” <https://github.com/osplanning-data-standards/GTFS-PLUS>.
- 40 29. “Dyno-demand Data Standard,” <https://github.com/osplanning-data-standards/dyno-demand>.