

# PondiõnsTracker: A framework based on GTFS-RT to identify delays and estimate arrivals dynamically in Public Transportation Network

Pedro Pongelupe Lopes

Programa de Pós-Graduação em Informática

December 04 2023



PUC Minas

# Contents

- 1 Introduction
- 2 Theoretical Reference
- 3 PondiônsTracker
- 4 Results
- 5 Conclusion



# Introduction

## Motivation

- Public Transportation Network
- Smart cities
- GTFS and GTFS-RT specifications



# Motivation

## GTFS-RT Matching Identifiers Issue

To work with GTFS-RT, it is **required** to track vehicles in real-time. But, in some cases, it is not easy to match the identifier between a real-time record and the GTFS static data. In Rome in 2016, this issue was reported by Raghothama et al. (2016). We still face this issue in Belo Horizonte in 2023.



# Objectives

## Main Objective

- Proposing and validating PondiônsTracker

## Specific Objectives

- Collecting data from the real-time API and combining with the GTFS
- Understanding if Belo Horizonte's delays are spatial and temporal dependent by analyzing delays among bus stops
- Comparing the arrival times defined at the GTFS with the arrival times generated by *PondiônsTracker*.



# Theoretical Reference

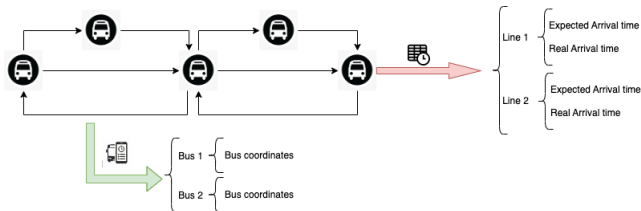
## Main Ideas

- Smart Cities
- Urban Computing
- Human Mobility
- GTFS and GTFS-RT
- Graphs and Complex Network



# Public Transportation Network as a Complex Network

$$G^I = (V_g, E_g, X_g, A_g)$$



$G^I$ : Graph  $G$  at a given time  $I$

$V_g$ : Bus stops

$E_g$ : Routes connecting two bus stops

$X_g$ : Additional information about bus stops ( $V_g$ )

$A_g$ : Additional information about routes connecting two bus stops ( $E_g$ )

# PondiônsTracker

## Overview

*PondiônsTracker*<sup>a</sup> is a framework to enrich GTFS data with real-time data. The name *PondiônsTracker* is a small gag from the sonority of the expression *bus stop* when pronounced in Portuguese with the accent from Minas Gerais.

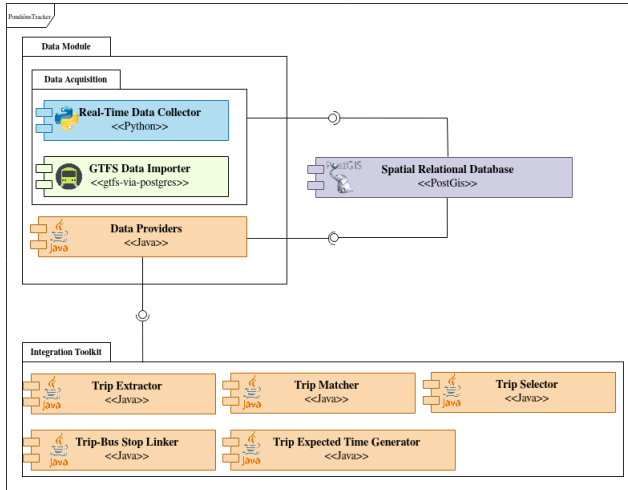
---

<sup>a</sup>Available at <https://github.com/Pongelupe/PondionsTracker/>

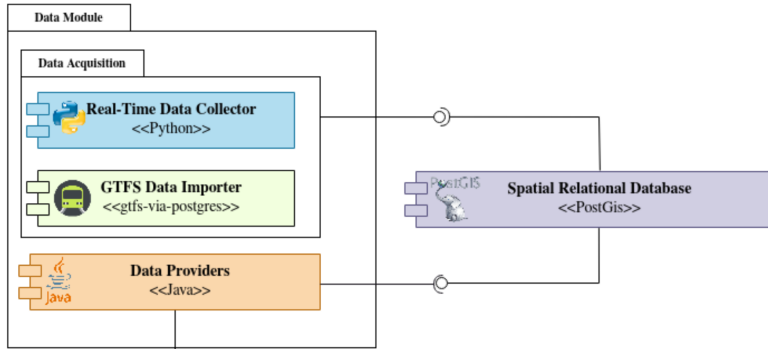




# PondiõnsTracker's Architecture



# Data Module Overview



# Data Providers

```
1  <dependency>
2    <groupId>br.pondionstracker</groupId>
3    <artifactId>data-module</artifactId>
4    <version>1.0.0</version>
5  </dependency>
```

Figura: *DataModule*'s maven dependency



# Integration Module

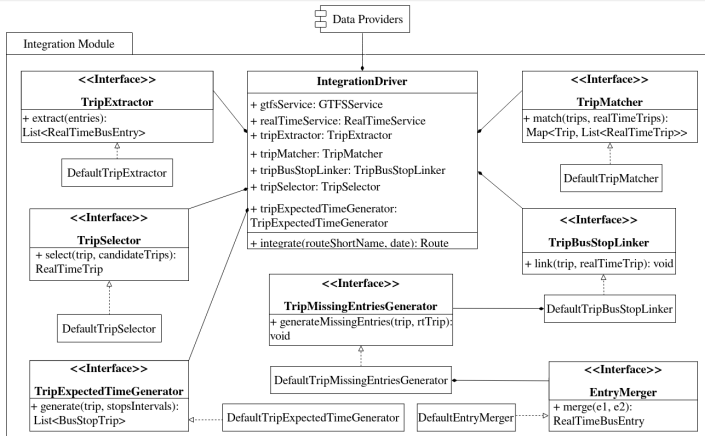


Figura: Integration Module Class Diagram

# Integration Driver

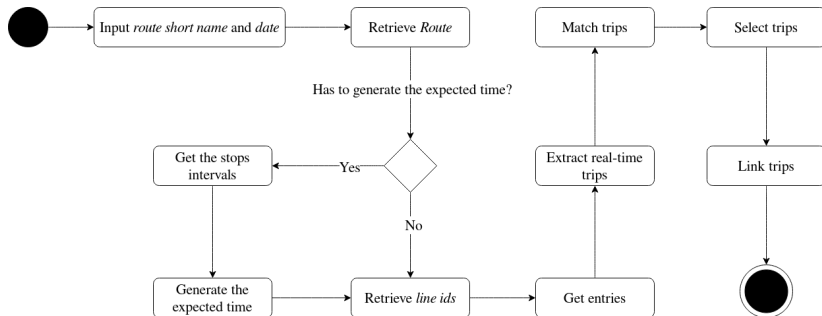


Figura: Integration Driver Activity Diagram

## Integration Driver - 1st Step

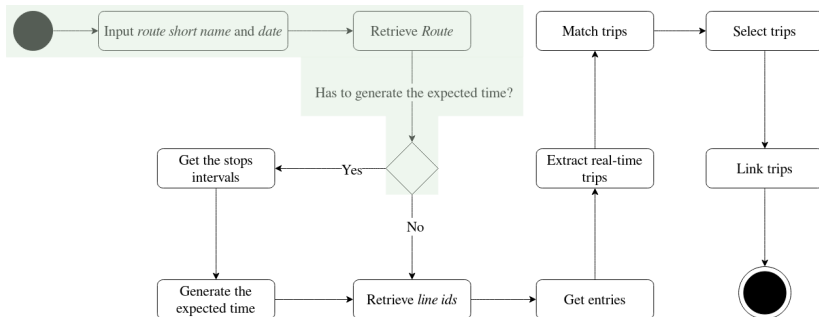


Figura: Integration Driver Activity Diagram

# Has to generate the expected time?

`arrival_time`

Time

Conditionally  
required

Arrival time at a specific stop for a specific trip on a route. If there are not separate times for arrival and departure at a stop, enter the same value for `arrival_time` and `departure_time`. For times occurring after midnight on the service day, enter the time as a value greater than 24:00:00 in HH:MM:SS local time for the day on which the trip schedule begins.

Scheduled stops where the vehicle strictly adheres to the specified arrival and departure times are timepoints. If this stop is not a timepoint, it is recommended to provide an estimated or interpolated time. If this is not available, `arrival_time` can be left empty. Further, indicate that interpolated times are provided with `timepoint=0`. If interpolated times are indicated with `timepoint=0`, then time points must be indicated with `timepoint=1`. Provide arrival times for all stops that are time points. An arrival time must be specified for the first and the last stop in a trip.

Figura: `arrival_time` definition from `stop_times.txt`



PUC Minas

## Integration Driver - 2nd Step\*

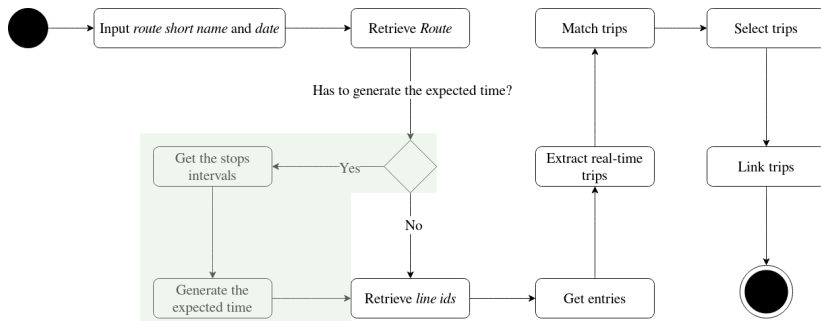


Figura: Integration Driver Activity Diagram



## Integration Driver - 3rd Step

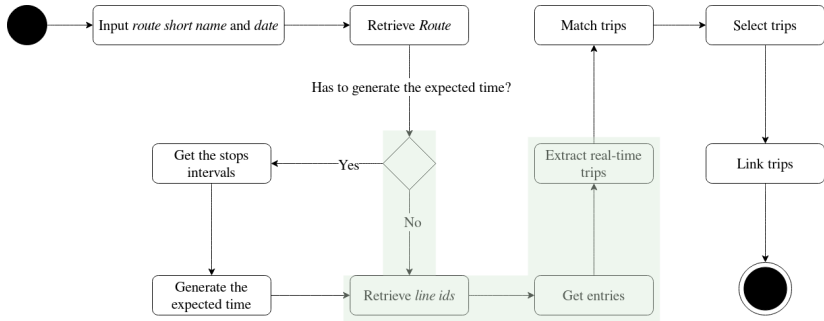


Figura: Integration Driver Activity Diagram

## Integration Driver - 4th Step

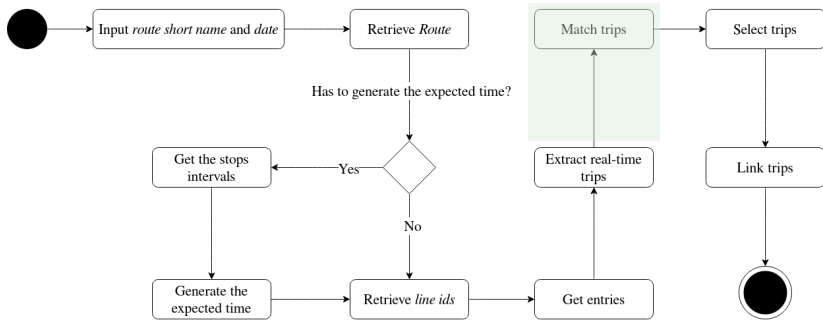


Figura: Integration Driver Activity Diagram

## Integration Driver - 5th Step

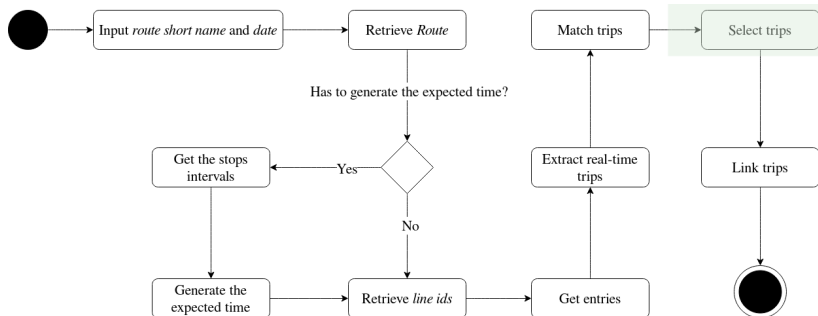


Figura: Integration Driver Activity Diagram

## Integration Driver - 6th Step

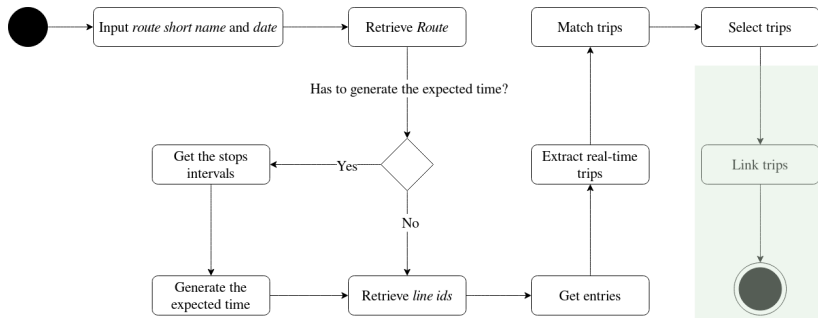


Figura: Integration Driver Activity Diagram

# Integration Module

```
1 <dependency>
2   <groupId>br.pondionstracker</groupId>
3   <artifactId>integration-module</artifactId>
4   <version>1.0.0</version>
5 </dependency>
```

Figura: *IntegrationModule*'s maven dependency



# PondiônsTracker-BH

## Overview

*PondiônsTracker-BH<sup>a</sup>* is our *PondiônsTracker*'s specialization to deal with Belo Horizonte's Network. So, we have implemented our own *Real-Time Data collector*, and we have overwritten a method from the *RealTimeService* from the *DataProviders*.

---

<sup>a</sup>Available at <https://github.com/Pongelupe/PondionsTracker-BH>



## Belo Horizonte's RealTimeService

### *BHRealTimeService - getIdsLineByRouteId*

This happens due to a **one-to-many** relationship between the GTFS and the real-time data.

- *BHTrans* → *GTFS*
- *Transfacil* → *Traffic API*



# Workload Overview

## Workload

- Data collected for 11 days straight in August 2023
- 30 Gigabytes

Date	Day-of-Week	Entries
29-07-23	Saturday	22,319,765
30-07-23	Sunday	22,635,117
31-07-23	Monday	22,583,380
01-08-23	Tuesday	22,432,739
02-08-23	Wednesday	21,970,073
03-08-23	Thursday	22,050,579
04-08-23	Friday	22,402,865
05-08-23	Saturday	22,642,955
06-08-23	Sunday	22,786,254
07-08-23	Monday	22,109,606
08-08-23	Tuesday	22,405,222
<b>Total</b>	-	<b>246,338,555</b>





# Schedule Analysis

## *Schedule-Filled Percentage*

*Schedule-Filled Percentage* = **Matched Trips / Scheduled Trips**

- **Total:**  $156,628 / 205,884 = 76.08\%$
- **Weekdays:**  $118,559 / 159,418 = 74.37\%$
- **Saturdays:**  $22,796 / 28,200 = 80.84\%$
- **Sundays:**  $15,273 / 18,266 = 83.61\%$



## Schedule Analysis - Schedule Deviations

### Schedule Deviations

Regarding the real-time API, there are collected trips which were not defined at Belo Horizonte's GTFS.

#### *82 - Estação São Gabriel / Savassi Via Hospitais*

On Sundays, the GTFS does not schedule any trip for route 82, but the API provided entries regarding this route twice during the period observed.

# Delay Analysis

## Delay Notation

- **Delay:**  $\geq 1$  minute after
- **Ahead-of-Schedule:**  $\geq 1$  minute before
- **On time:**  $\leq 59$  seconds after OR  $\leq 59$  seconds before

	Weekday	Saturday	Sunday
Total trips matched	118,559	22,796	15,273
Trips entirely out of schedule	60,244	10,899	7,148
Trips with departure or arrival on time	39,403	8,731	5,988
Trips with departure and arrival on time	324	95	56
Trips entirely on time	1	2	1

**Figura:** Delays detailed in whole Public Transportation Network scale



## Delay Analysis

### *331 - Estação Barreiro/Conjunto Antonio Teixeira Dias Via Upa*

Has 32 bus stops, representing a length of almost 9 kilometers.

- ① Jul. 29 15:30:00 - 15:56:27 → Jul. 29 15:30:03 - 15:57:03
- ② Jul. 30 08:20:00 - 08:46:27 → Jul. 30 08:20:31 - 08:46:15
- ③ Aug. 04 05:40:00 - 06:06:27 → Aug. 04 05:40:30 - 06:06:00
- ④ Aug. 05 17:10:00 - 17:36:27 → Aug. 05 17:10:45 - 17:36:49

# Delay Analysis

## Distribution of each status over the network

- **Delay:** 89.8%
- **Ahead-of-Schedule:** 6.9%
- **On time:** 3.3%

## Attention!

The predominance of *DELAYED* in the Public Transportation Network **does not imply** that the network is not working nor completely stopped!

# Delay Analysis

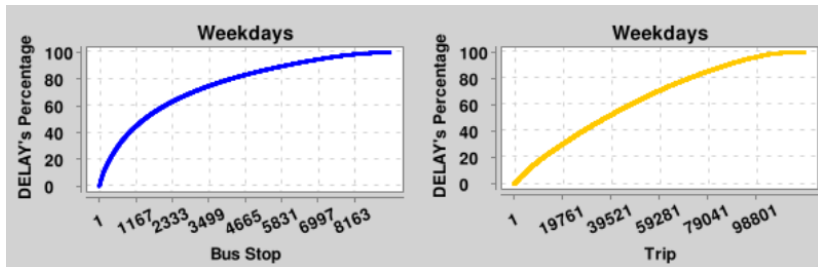


Figura: *DELAY*s Distribution: Bus Stop and Trip

# Delay Analysis

Figura: 300 Most Delayed Stops

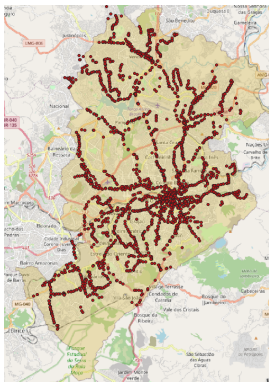
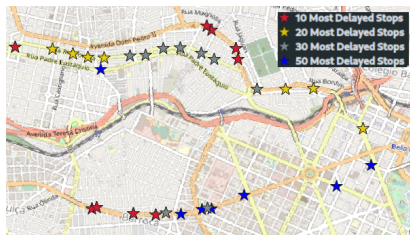


Figura: Fragment of the 50 Most Delayed Stops



# Delay Analysis

## Three Most Delayed Stops for Weekdays

- 1 #14793268 - *Avenida Dom Pedro II 1520* with 7,309 delays
- 2 #14791617 - *Avenida Amazonas 7309* with 7,009 delays
- 3 #14790997 - *Avenida Dom Pedro II 1980* with 6,692 delays

## Constants

- 1 *Global Ahead Average*: 13.42 minutes
- 2 *Global Delay Average*: 20.49 minutes



## Delay Analysis

### Stops #14793268 and #14790997

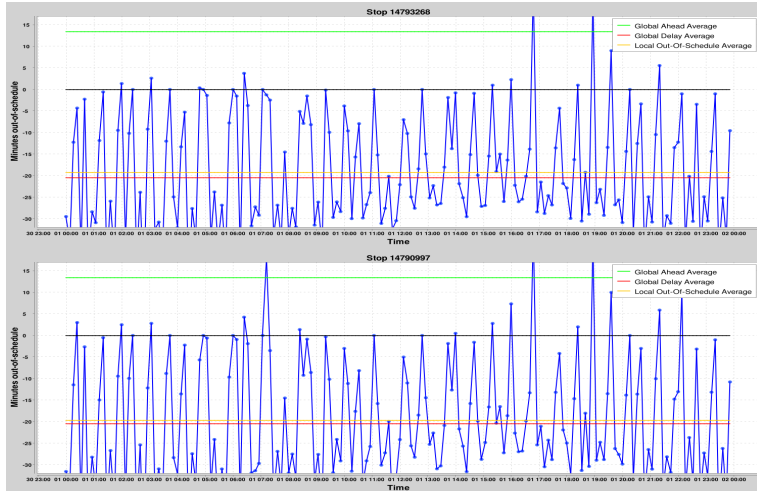
The stops #14793268 and #14790997 are the first and third most delayed in the Public Transportation Network, respectively. Also, these stops are **462** meters from each other on the same avenue, *Avenida Pedro II*, and share **2,590** common trips, so they are spatially related.

### Local Out-Of-Schedule Average

- #14793268: 19.29 minutes
- #14790997: 19.68 minutes



# Delay Analysis



# Comparison Between Generated and Real Data

## Overview

The previous analysis was only possible because Belo Horizonte's GTFS defines the expected time for all bus stops on every trip. The *Trip Expected Time Generator* generates the expected times when missing, so, we executed this component with Belo Horizonte's data and compared the expected times generated with those defined at the GTFS.



## Comparison Between Generated and Real Data

### Trip entirely out of schedule

- GTFS: 78,291
- Generated: 75,073
- **Diff: 3,218 (4.29%)**

### Trip entirely on time

- GTFS: 4
- Generated: 0
- **Diff: 4 (100%)**

### Trips with departure or arrival on time

- GTFS: 54,122
- Generated: 54,271
- **Diff: 149 (0.27%)**

### Trips with departure and arrival on time

- GTFS: 475
- Generated: 596
- **Diff: 121 (25.47%)**

## Comparison Between Generated and Real Data

		GTFS	Generated
Weekday	<i>ON_TIME</i>	3.3%	3.2%
	<i>AHEAD_OF_SCHEDULE</i>	6.9%	17.8%
	<i>DELAYED</i>	89.8%	79.0%
Saturday	<i>ON_TIME</i>	3.9%	3.5%
	<i>AHEAD_OF_SCHEDULE</i>	6.5%	18.4%
	<i>DELAYED</i>	89.6%	78.1%
Sunday	<i>ON_TIME</i>	4.4%	3.9%
	<i>AHEAD_OF_SCHEDULE</i>	5.4%	18.5%
	<i>DELAYED</i>	90.2%	77.6%

# Comparison Between Generated and Real Data

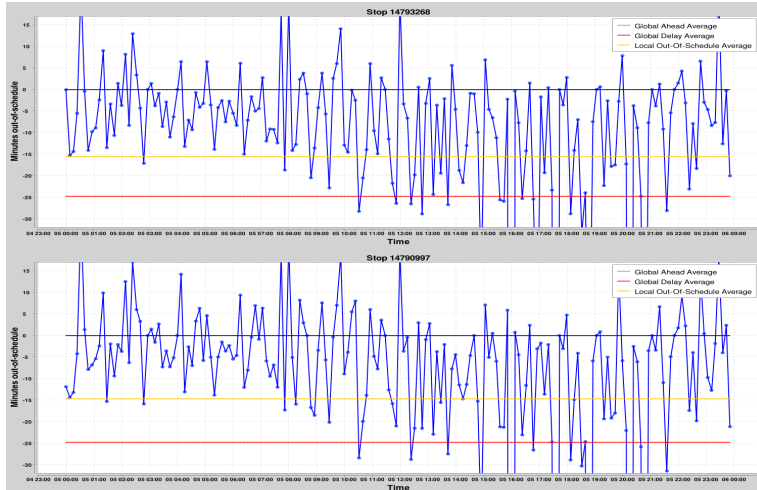
## Global Averages

- *Global Ahead Average*
  - GTFS: 13.42 minutes
  - Generated: 38.57 minutes
  - **Diff:** 25.15 minutes
- *Global Delay Average*
  - GTFS: 20.49 minutes
  - Generated: 24.75 minutes
  - **Diff:** 4.26 minutes

## Local Out-Of-Schedule Average

- #14793268
  - GTFS: 19.29 minutes
  - Generated: 15.54 minutes
  - **Diff:** 3.75 minutes
- #14790997
  - GTFS: 19.68 minutes
  - Generated: 14.68 minutes
  - **Diff:** 5 minutes

# Comparison Between Generated and Real Data



# Limitations

## Limitations

The *Real-Time Data Collector* is the most fragile component due to the third-party real-time traffic API interface.

- Size and quality of the data
- Scheduled routes with no entries reported
  - 1 720 - *Circular Saúde MG20* missed 175 trips
  - 2 912 - *Conjunto Taquaril/Praça Che Guevara* missed 210 trips



# Conclusion

## Concluding Remarks

- Comparison between the expected schedule with the actual schedule
- Delays in Belo Horizonte follow a *log-normal* distribution
- Analysis using data generated with the *Trip Expected Time Generator*
- *PondiônsTracker* as a viable option when GTFS-RT is unavailable

# Conclusion

## Future Work

- Further explore Belo Horizonte Public Transportation Network using deep learning for graphs approaches
- Reproduce Belo Horizonte's results with other cities
- Explore the delays analysis combining temporal and spatial dimensions

Jayanth Raghothama, Vinutha Magal Shreenath, and Sebastiaan Meijer. 2016. Analytics on Public Transport Delays with Spatial Big Data. In *Proceedings of the 5th ACM SIGSPATIAL International Workshop on Analytics for Big Geospatial Data* (Burlingame, California) (*BigSpatial '16*). Association for Computing Machinery, New York, NY, USA, 28–33.  
<https://doi.org/10.1145/3006386.3006387>

# Conclusion

Thanks!!



PUC Minas