# Transformation ML Framework
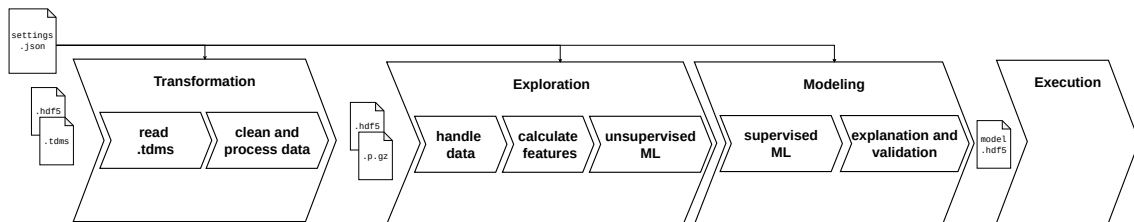
## on basis of the XBox2 Data Set

Lorenz Fischl

# Introduction

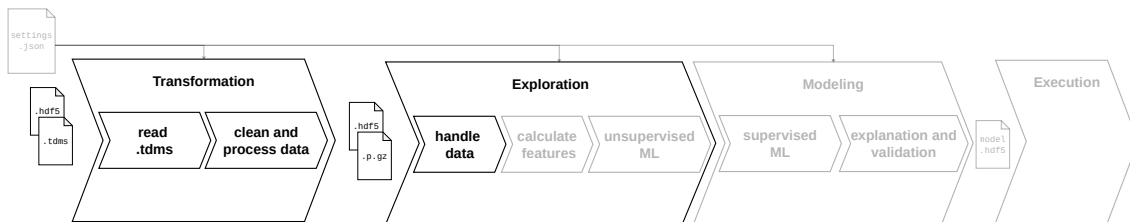# Introduction

# Table of Contents

# Table of Contents

# Table of Contents

# XBox2 Data

# XBox2 Data

Every 20$ms$ a pulse is sent into the RF cavity for particle acceleration.

# XBox2 Data

Every 20*ms* a pulse is sent into the RF cavity for particle acceleration.
Sometimes an arc forms. Those events are called breakdown.

# XBox2 Data

Every 20*ms* a pulse is sent into the RF cavity for particle acceleration. Sometimes an arc forms. Those events are called breakdown.
vspace1cm



- ■ healthy pulse
- ■ breakdown pulse

} represents one pulse, that is one continuouse time series

data measured →

$\left[\texttt{double}^{3200}\right]^8$

$\left[\texttt{double}^{500}\right]^8$

# XBox2 EventData

# XBox2 EventData

A log group of one pulse is stored every minute.

# XBox2 EventData

A log group of one pulse is stored every minute.
When a breakdown happens the corresponding log group + the two
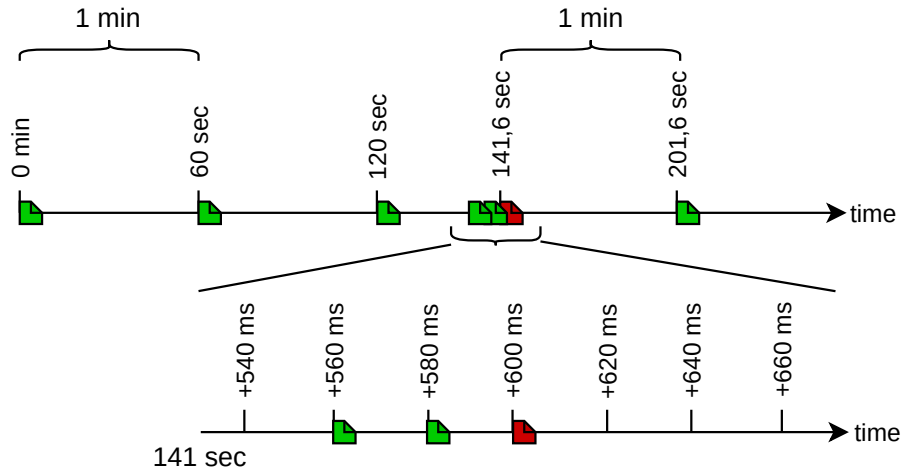prior log groups are stored.

# XBox2 EventData

A log group of one pulse is stored every minute.
When a breakdown happens the corresponding log group + the two prior log groups are stored.

# XBox2 TrendData

# XBox2 TrendData

35 values about the environmental conditions (that don't change rapidly) are stored roughly every 1,5 sec.

# XBox2 TrendData

35 values about the environmental conditions (that don't change rapidly) are stored roughly every 1,5 sec.
All TrendData of one day is stored in 1-2 groups.

# XBox2 TrendData

35 values about the environmental conditions (that don't change rapidly) are stored roughly every 1,5 sec.
All TrendData of one day is stored in 1-2 groups.



environmental data

# Table of Contents

# Table of Contents

# nptdms

# nptdms

- package `nptdms` can read and handle `.tdms` files

# nptdms

- package `nptdms` can read and handle `.tdms` files
- very slow (ex.: read of a 100MB file can take >30sec)

# nptdms

- package `nptdms` can read and handle `.tdms` files
- very slow (ex.: read of a 100MB file can take >30sec)
- very space inefficient (ex. TrendData: 20,5 GB in .tdms $\rightarrow$ 2.8 GB of data)

# Table of Contents

# Transformation

# Transformation

.tdms $\longrightarrow$ pd df/ dictionary $\longrightarrow$ pickle $\longrightarrow$ .gzip
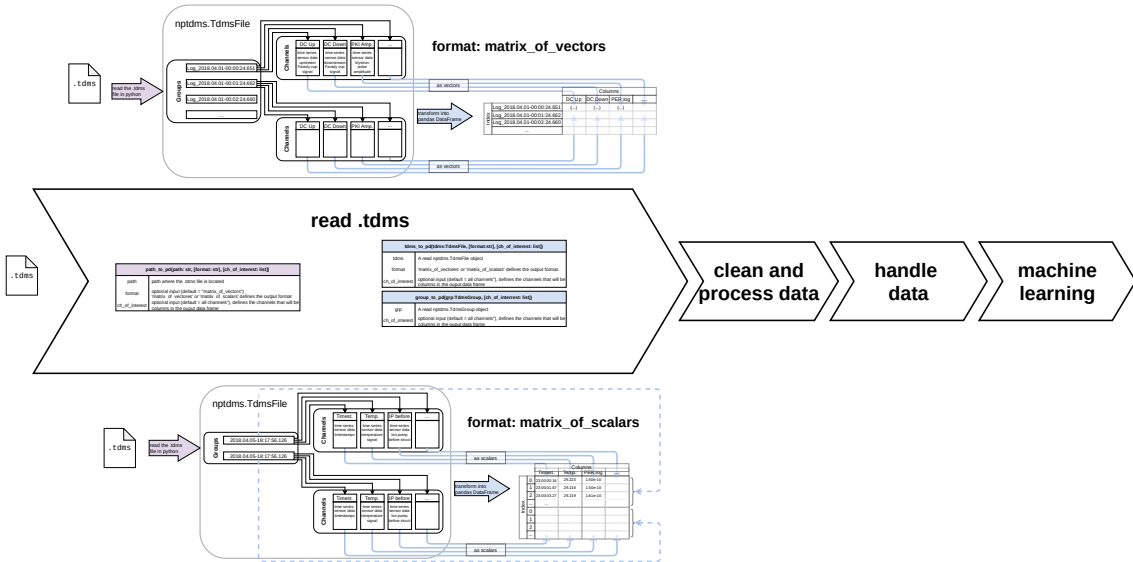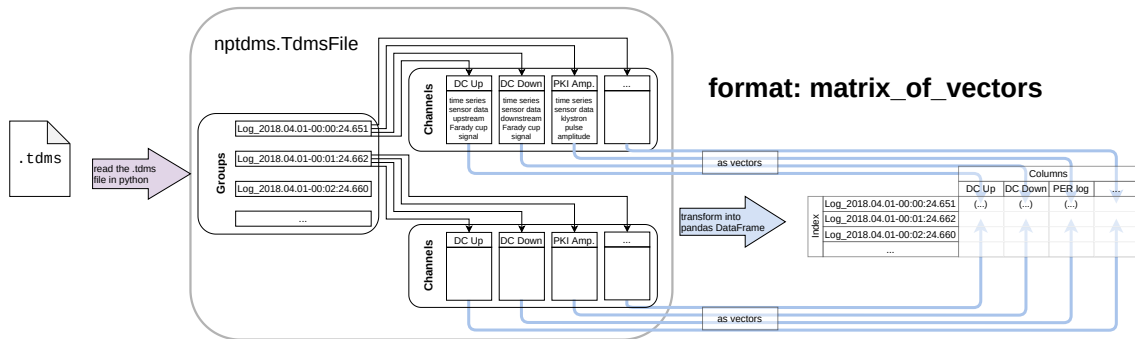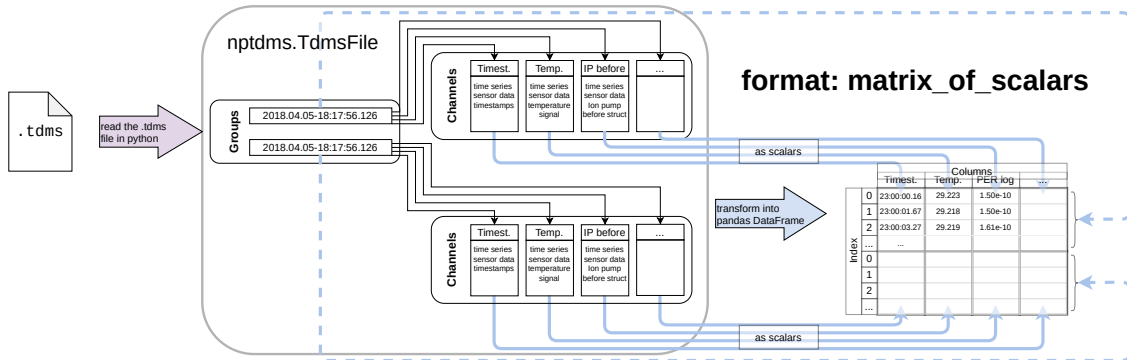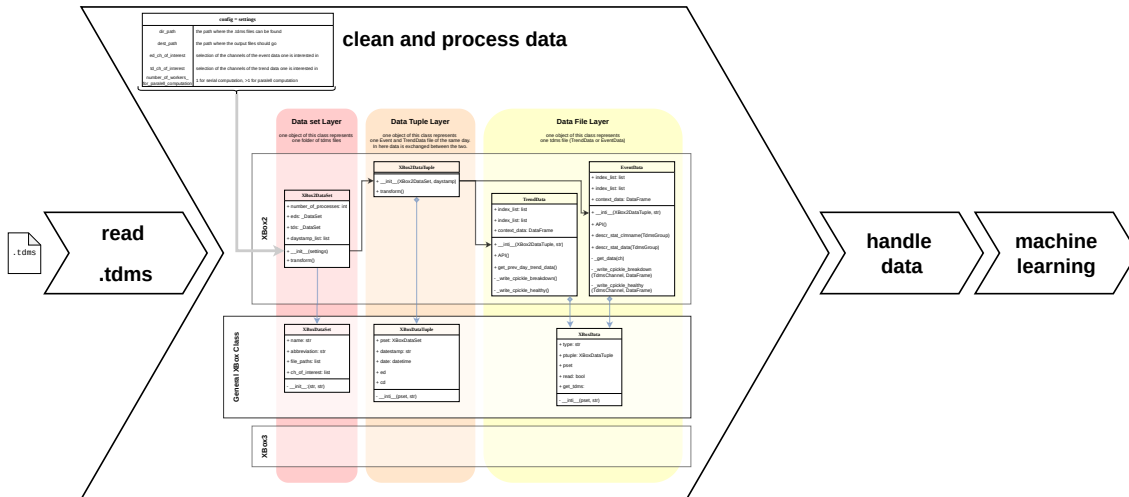
# Transformation: read.tdms

# Transformation: read.tdms

# Transformation: read.tdms

# Transformation: Clean and Process data with classes

Data set Layer
one object of this class represents
one folder of tdms files

Data Tuple Layer
one object of this class represents
one Event and TrendData file of the same day.
In here data is exchanged between the two.

Data File Layer
one object of this class represents
one tdms file (TrendData or EventData)

**XBox2**

**XBox2DataSet**

+ number_of_processes: int
+ eds: _DataSet
+ tds: _DataSet
+ daystamp_list: list

+ __init__(settings)
+ transform()

**XBox2DataTuple**

+ __init__(XBox2DataSet, daystamp)
+ transform()

**TrendData**

+ index_list: list
+ index_list: list
+ context_data: DataFrame
+ __inti__(XBox2DataTuple, str)
+ API()
+ get_prev_day_trend_data()
- _write_cpickle_breakdown()
- _write_cpickle_healthy()

**EventData**

+ index_list: list
+ index_list: list
+ context_data: DataFrame
+ __inti__(XBox2DataTuple, str)
+ API()
+ descr_stat_clmname(TdmsGroup)
+ descr_stat_data(TdmsGroup)
- _get_data(ch)
- _write_cpickle_breakdown
(TdmsChannel, DataFrame)
- _write_cpickle_healthy
(TdmsChannel, DataFrame)

**General XBox Class**

**XBoxDataSet**

+ name: str
+ abbreviation: str
+ file_paths: list
+ ch_of_interest: list
- __init__:(str, str)

**XBoxDataTuple**

+ pset: XBoxDataSet
+ datestamp: str
+ date: datetime
+ ed
+ cd
- __inti__(pset, str)

**XBoxData**

+ type: str
+ ptuple: XBoxDataTuple
+ pset
+ read: bool
+ get_tdms:
- __inti__(pset, str)

**XBox3**

# Summary

# Summary

- pandas DataFrame are easy to use in notebooks

# Summary

- `pandas` DataFrame are easy to use in notebooks
- pickle speeds up reading time

# Summary

- pandas DataFrame are easy to use in notebooks
- pickle speeds up reading time
- with compression takes up less space

# Summary

- pandas DataFrame are easy to use in notebooks
- pickle speeds up reading time
- with compression takes up less space


- version issues with pickle protocol

# Summary

- `pandas` DataFrame are easy to use in notebooks
- pickle speeds up reading time
- with compression takes up less space


- version issues with pickle protocol
- part of the process should not be in the Datahanlder instead of the Transformation

# Summary

- pandas DataFrame are easy to use in notebooks
- pickle speeds up reading time
- with compression takes up less space


- version issues with pickle protocol
- part of the process should not be in the Datahanlder instead of the Transformation
- EventData and TrenadData are stored differently

# Summary

- pandas DataFrame are easy to use in notebooks
- pickle speeds up reading time
- with compression takes up less space

- version issues with pickle protocol
- part of the process should not be in the Datahanlder instead of the Transformation
- EventData and TrenadData are stored differently
- channel properties are lost

# Summary

- pandas DataFrame are easy to use in notebooks
- pickle speeds up reading time
- with compression takes up less space

<br>

- version issues with pickle protocol
- part of the process should not be in the Datahanlder instead of the Transformation
- EventData and TrenadData are stored differently
- channel properties are lost
- data was changed in place in notebooks in retrospect

# Table of Contents

# Conclusion

# Conclusion

- I implemented a generic class for converting `.tdms` files into `pd.df+cpickle`

# Conclusion

- I implemented a generic class for converting `.tdms` files into `pd.df+cpickle`
- is there a better data format, maybe `.hdf5`?

# Conclusion

- I implemented a generic class for converting `.tdms` files into `pd.df+cpickle`
- is there a better data format, maybe `.hdf5`?

# Conclusion

- I implemented a generic class for converting `.tdms` files into `pd.df+cpickle`
- is there a better data format, maybe `.hdf5`?

|  | nptdms | pd.df+cpickle | | .hdf5 | |
|---|---|---|---|---|---|
|  |  | w/o zip | w zip | w/o zip | w/ zip |
| space (GByte) | 20.5GB | 2.8GB | 1GB | 2.8GB | 1GB |
| read (TD 1 channel) | $\sim$60min | 4sec | 12sec | 0.5 sec | |
| read (TD 3 channels) | $\sim$60min | 4sec | 12sec | 1 sec | |
| read (TD 15 channels) | $\sim$60min | 4sec | 12sec | 4 sec | |
| feature calc. (ED) | $>$15min | | 7sec | 8 sec | |

home.cern