Ceatech to industry





STREAMER A POWERFUL AND OPEN-SOURCE FRAMEWORK FOR CONTINUOUS LEARNING

IN DATA STREAMS



Sandra GARCIA RODRIGUEZ







list

Paradigm shift

2 consequences

- Get a model while data is still incoming
- Non-stationarity





• Definition

- Infinite streams of data integrated from both live and historical sources
- Stream: $[d_0, d_1 \dots d_n]$



- They contain entities [and labels] <e,L>
- Each entity *e* is linked to *1..n* instances (pieces of information)

	Batch processing	Stream processing
Data scope	Queries or processing over all or most of the data in the dataset.	Queries or processing over data within a rolling time window, or on just the most recent data record.
Data size	Large batches of data.	Individual records or micro batches consisting of a few records.
Performance	Latencies in minutes to hours.	Requires latency in the order of seconds or milliseconds.
Analyses	Complex analytics.	Simple response functions, aggregates, and rolling metrics.

- **DS processing tools** (as STREAMER)
 - algorithms applied to DS that can learn incrementally



list



WHAT IS STREAMER?









<u>Data Streams</u>: data that arrive continuously but cannot be stored, it needs to be processed in real time.

- □ What is it? framework for integrating and testing machine learning algorithms in realistic streaming operational contexts
 - □ in multiple programming languages
 - □ some already available
- □ Users can abstract from streams implementation (transparent for them) to only focus on the algorithms



Implementation of the whole pipe (or just part of it) of data stream ingestion, processing, knowledge extraction and visualization

Features:

- open source
- flexible & scalable
- □ cross-platform
- logs
- visualization

eatech objectives

Control - machine learning experimenter

- Allows users to easily integrate and test machine learning algorithms into realistic streaming operational context
- □ User must be able to define the context and learning configuration
- Provides several functionalities (streaming scenario simulation, train, test, evaluation, models inference, monitoring...)

Easiness

- Easy installation and deployment
- Visualization capabilities

Reactivity

Process and incrementally learn from continuous data streams





Robust and secure deployment

Ceatech STREAM LEARNING FUNCTIONALITIES



Ceatech STREAMER IN DETAILS

list

□ The code is published in open-source on GitHub

Official Website: <u>https://streamer-framework.github.io/</u>

STREAMER

a Powerful and Open-Source Framework for Continuous Learning in Data Streams



What is STREAMER?

STREAMER is a modern framework that helps scientists to easily integrate and test machine learning algorithms into realistic streaming operational contexts.



STREAMER MORE IN DETAIL



www.cea.fr







- License (use and distribution): terms by GNU GENERAL PUBLIC LICENSE Version 3 (<u>https://www.gnu.org/licenses/gpl-3.0.html</u>)
- Easy **installation**:
 - Docker environment
 - □ Manual choice (advanced mode)



Ceatech STREAMER ARCHITECTURE

□ STREAMER consists of a layered architecture (independent modules):

- M1 Data Stream Ingestion
 - □ Offline data are sent by blocks of N records and periodically every t seconds
 - Online data are sent according its timestamps, sending can be scaled, advanced, or delayed in time.
- □ M2 DSPE Processor
 - Drives the streaming pipeline from M1 to M5
- □ M3 Learning API/ API adapter
 - Supports batch training and online training (learning update)
- □ M4 Applications & Evaluation
 - Pre/post processing functionalities & evaluation metrics (accuracy, sensitivity, precision, etc.)
- □ M5 Data Visualization
 - □ The model's performance can be monitored in a dedicated Kibana dashboard.











STREAMER DISSEMINATION







www.cea.fr





- Garcia-Rodriguez, Sandra, Mohammad Alshaer, and Cedric Gouy-Pailler.
 "STREAMER: A Powerful Framework for Continuous Learning in Data Streams." Proceedings of the 29th ACM International Conference on Information & Knowledge Management. 2020.
 <u>https://dl.acm.org/doi/10.1145/3340531.3417427</u>
- Alshaer, Mohammad, Sandra Garcia-Rodriguez, and Cedric Gouy-Pailler.
 "Detecting Anomalies from Streaming Time Series using Matrix Profile and Shapelets Learning." 2020 IEEE 32nd International Conference on Tools with Artificial Intelligence (ICTAI). IEEE, 2020.
 <u>https://www.computer.org/csdl/proceedingsarticle/ictai/2020/922800a376/1pP3tGGtcn6</u>

Ceatech OPEN PUBLIC

li/1

- CEA Tech Linkedin = https://www.linkedin.com/feed/update/urn:li:activity:6787752698235224064/
- CEA Tech News (EN) = <u>https://www.cea-tech.fr/cea-tech/english/Pages/2021/streamer-simulates-massive-data-streams.aspx</u>
- CEA Tech News (FR) = <u>https://www.cea-tech.fr/cea-tech/Pages/2021/streamer-simule-des-flux-de-donnees-massives.aspx</u>
- Linkedin (FR/EN) = https://www.linkedin.com/feed/update/urn:li:activity:6787752698235224064/
- Twitter (FR) = <u>https://twitter.com/CEA_List/status/1377182719638003715?s=20</u>
- Annual report (FR) = <u>http://www-list.cea.fr/images/stories/decouvrir-le-list/qui-sommes-nous/rapport-dactivite/Rapport_activite_CEA_LIST_2020.pdf</u>
- CEA LIST News (FR) = <u>http://www-list.cea.fr/medias/toute-l-actualite/2021/496-30-mars-</u> 2021-streamer-une-plateforme-logicielle-au-service-de-l-apprentissage-automatique-sur-les-<u>flux-de-donnees</u>
- CEA LIST News (EN) = <u>http://www-list.cea.fr/en/media/news/2021/496-march-30-2021-</u> <u>streamer-a-software-platform-for-machine-learning-for-data-streams</u>
- Paris Saclay University (EN) : <u>https://www.universite-paris-saclay.fr/en/news/streamer-software-platform-machine-learning-data-streams</u>
- DataIA (FR/EN) : <u>https://www.actuia.com/actualite/streamer-un-programme-permettant-dintegrer-et-de-tester-facilement-des-algorithmes-de-machine-learning/</u>



STREAMER LIFE DEMO



www.cea.fr





- Objective: electrocardiogram (ECG) anomaly detection in time series stream data with no prior labels
 - anomalies may vary in time

MIT-BIH database

- □ 48 ambulatory ECG 30 mins recordings
- every heartbeat is represented by a different sequence
- □ Anomaly: Ventricular ectopic beat class

Unsupervised learning

- Matrix Profile algorithm for learning the shapelets (features) that represent anomalies over unsupervised data streams
- Based on them, we do the detection of the anomalies by performing a similarity function
- Model is updated after every bunch of data received

Evaluation: labels of MIT-BIH are only used for metrics calculation

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

 $Sensitivity = \frac{TP}{TP + FN}$ where TP means true positive, TN true negative, FP false positive and FN false negative.

FROM RESEARCH TO INDUSTRY



list



Demonstration video is found at the following link:

https://www.youtube.com/watch?v=UHx1YenTF1Y&feature=emb_title

STREAMER can run several scenarios in parallel
 Different configuration files for each UC

- Network intrusion detection (KDD)
 - □ Classification (23 classes)
 - □ Algorithm: Hoeffding Tree Classifier
 - □ Metrics: accuracy, F1 score, precision & sensitivity
- Demand forecasting (Multivariate time series)
 - Regression
 - Algorithm: XGBoost
 - □ Metrics: MSE, RMSE, MAE, R2



10:35:00 10:36:00 10:37:00 10:38:00 10:39:00 10:40:00 10:41:00 10:42:00 10:43:00 10:44:00 10:45:00 10:45:00 10:46:00 10:47:00 10:48:00 **per 10 seconds**

0 10:34:00





STREAMER

a Powerful and Open-Source Framework for Continuous Learning in Data Streams

Commissariat à l'énergie atomique et aux énergies alternatives Institut Carnot CEA LIST CEA Saclay – DIGITEO LABS | 91191 Gif-sur-Yvette Cedex Phone. +33 (0)1 69 08 88 55 | Fax. +33 (0)1 69 08 89 65 Direction: CEA TECH

Department: DCSI Laboratory: LI3A

Etablissement public à caractère industriel et commercial RCS Paris B 775 685 019