



FORTISSIMO



# Easy use of Distributed TensorFlow Training on supercomputing facilities

Gonzalo Ferro, CESGA  
gferro@cesga.es

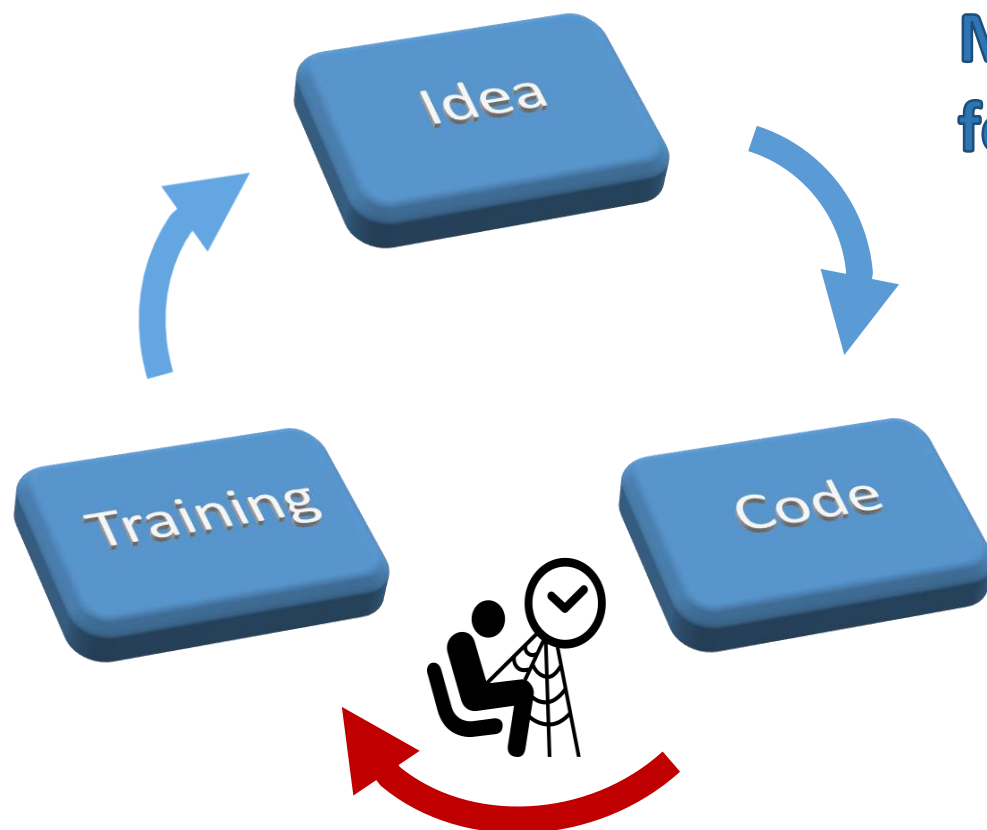
IBERGRID



IBERGRID 2018: Towards the European Open Science Cloud – EOSC.  
11th - 12th October. Lisbon, ISCTE - University Institute of Lisbon (ISCTE-IUL)



# Machine Learning design cycle



**Machine Learning (ML) is a powerful tool for science, industry and other sectors.**

**Performance of ML algorithms is improved by training them using large datasets.**

**HPC can help engineers to boost their algorithms.**

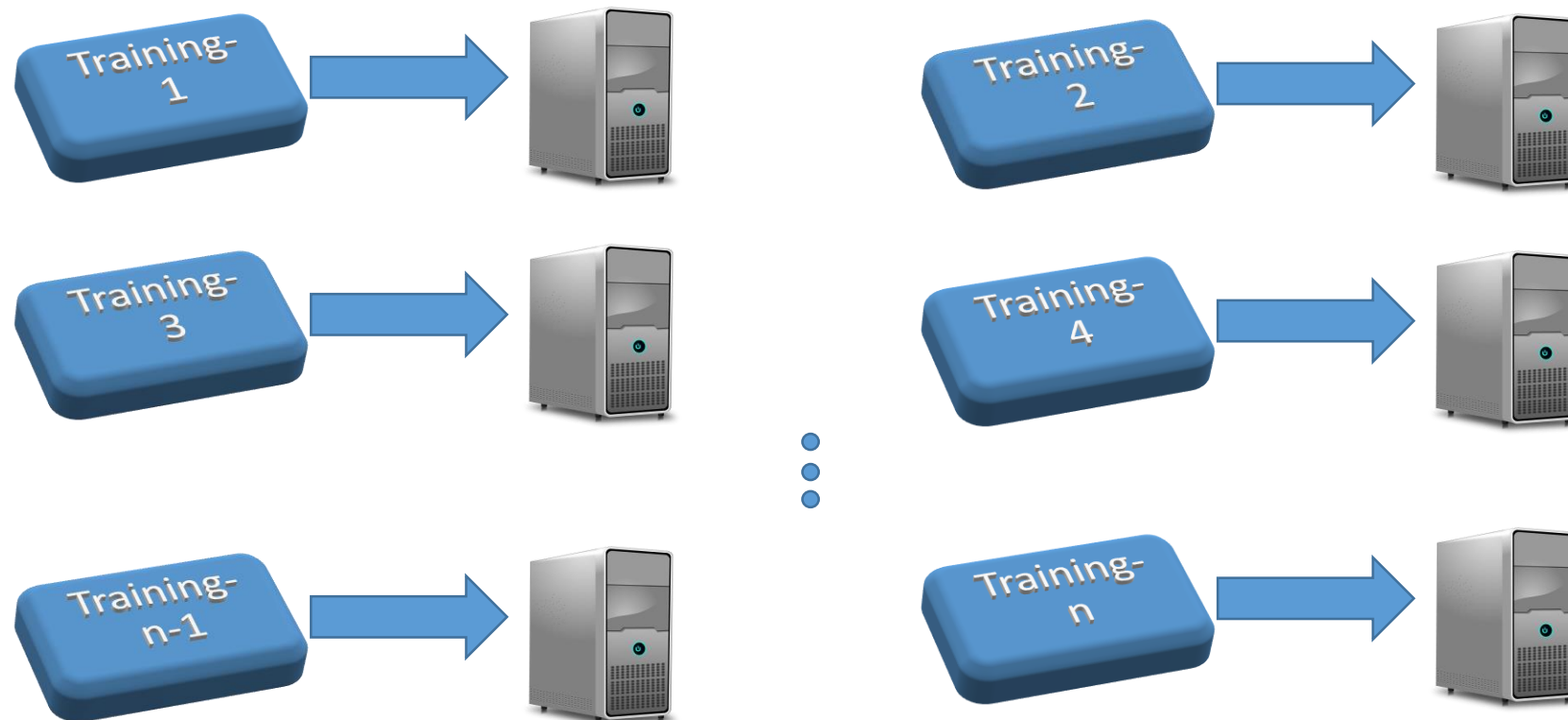
# How exploit HPC for ML training?

- Simultaneous Training.
- Parallel Distributed Training.
- Simultaneous + Parallel Distributed.



# How exploit HPC for ML training?

## Simultaneous Training



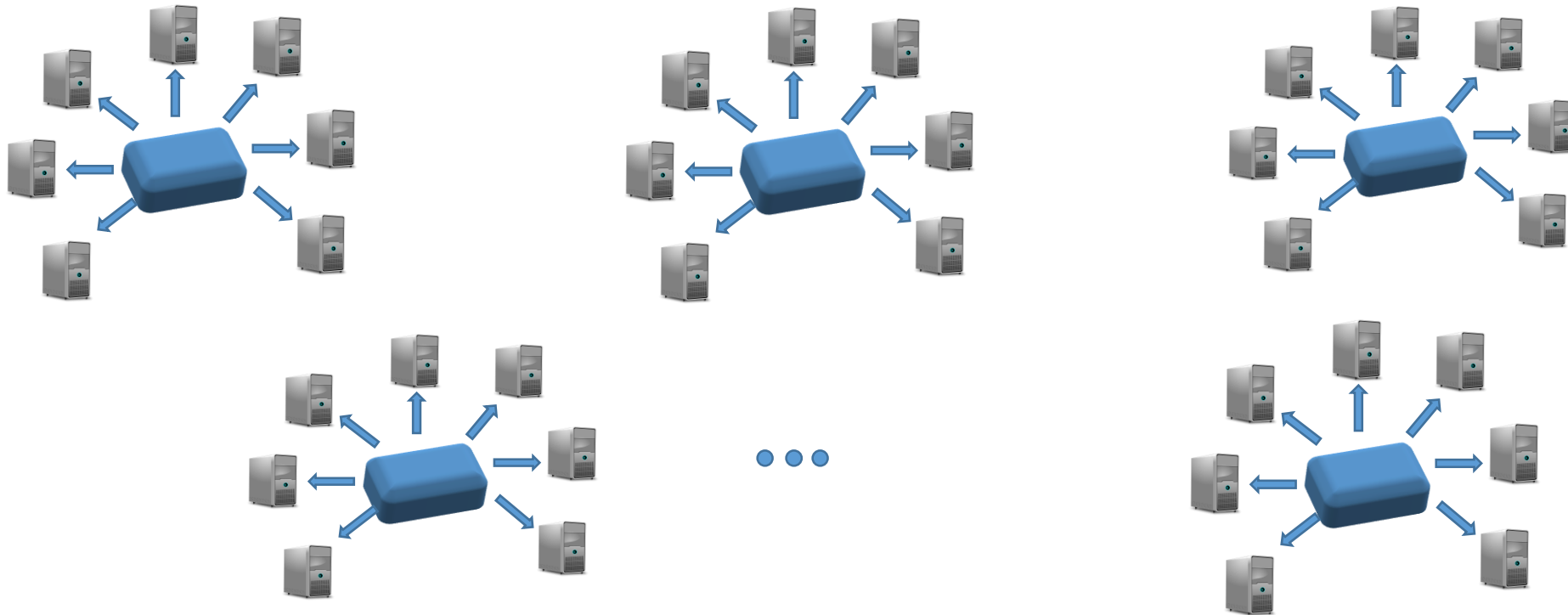
# How exploit HPC for ML training?

## Parallel Distributed Training



# How exploit HPC for ML training?

## Simultaneous + Parallel Distributed



# How exploit HPC for ML training?

- Simultaneous Training
- **Parallel Distributed Training**
- Simultaneous + Parallel Distributed



# Distributed Training API



**TensorFlow (TF)**

- **A Machine Learning API developed by Google.**
- **One of the most widely of the tools used for developing and training of deep learning models.**
- **TF allows users to implement distributed computing capabilities in their training in an easy way.**





# Distributed Training API

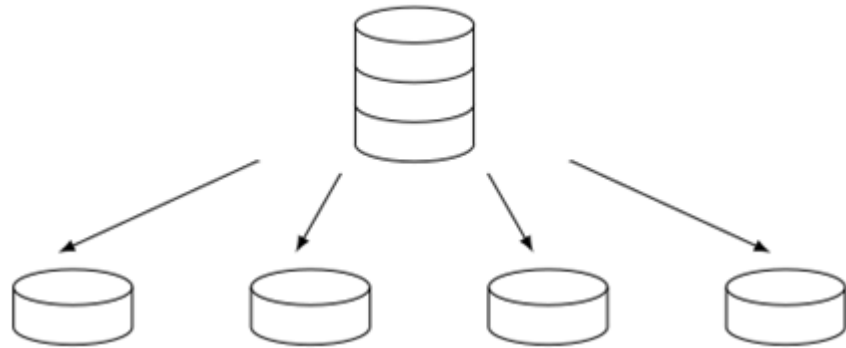


TensorFlow (TF)

- A Machine Learning API developed by Google.
- One of the most widely of the tools used for developing and training of deep learning models.
- TF allows users to implement **distributed computing capabilities** in their training in an easy way.



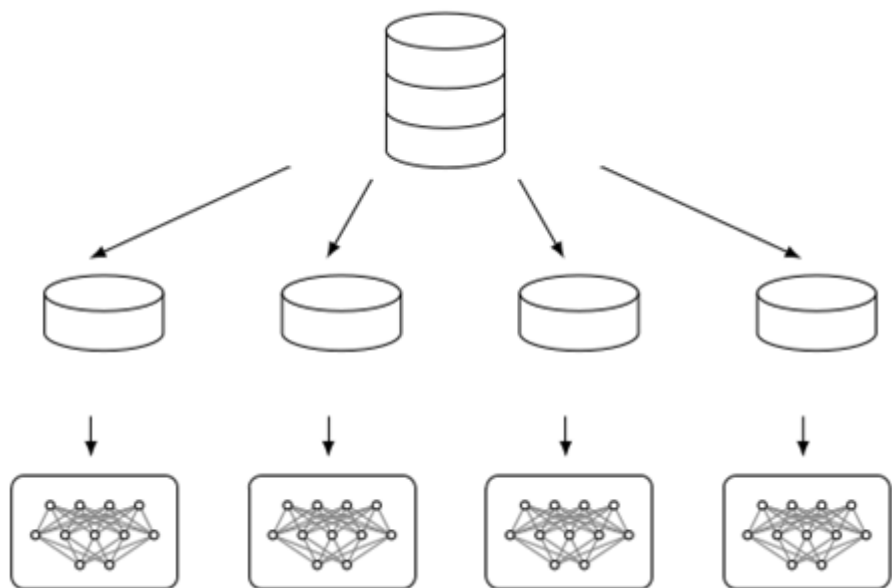
# Distributed TF: Data Parallelism.



**Splitting of Training Dataset.**



# Distributed TF: Data Parallelism.



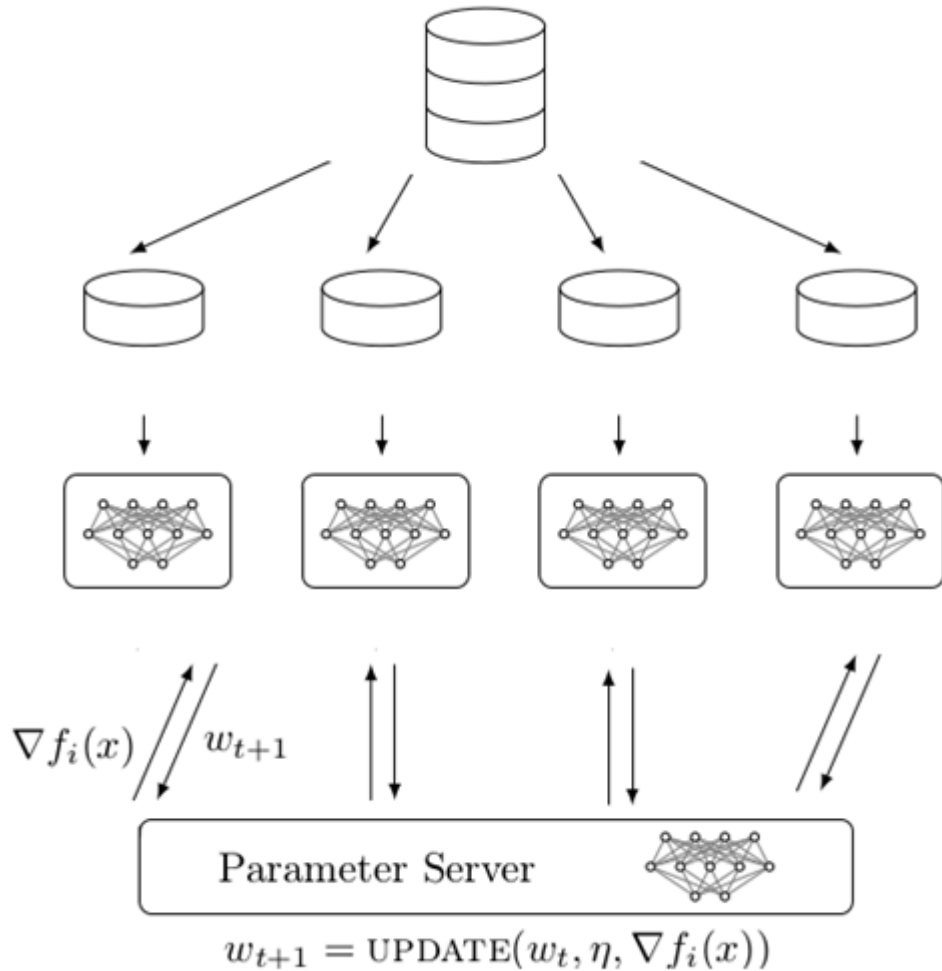
**Splitting of Training Dataset.**

**Workers:**

- Copy of the computational graph.
- Calculate gradients over their correspondent part of dataset.



# Distributed TF: Data Parallelism.



**Splitting of Training Dataset.**

**Workers:**

- Copy of the computational graph.
- Calculate gradients over their correspondent part of dataset.



**Parameter Servers:**

- Store weights and bias of the model.
- Responsible for the aggregation of gradients calculated by Workers.



# Distributed TF: Queue system issues.

Distributed TF is based on communication protocol called gRPC.



# Distributed TF: Queue system issues.

Distributed TF is based on communication protocol called gRPC.

CESGA Finis Terrae II (FT2) uses Slurm for Resource Management.



# Distributed TF: Queue system issues.

Distributed TF is based on communication protocol called gRPC.

CESGA Finis Terrae II (FT2) uses Slurm for Resource Management.

- Detected issues when deploying Distributed TF on FT2:



# Distributed TF: Queue system issues.

Distributed TF is based on communication protocol called gRPC.

CESGA Finis Terrae II (FT2) uses Slurm for Resource Management.

- Detected issues when deploying Distributed TF on FT2:
  1. Distributed TF needs **Network Addresses (IP:Port)** in advance.





# Distributed TF: Queue system issues.

Distributed TF is based on communication protocol called gRPC.

CESGA Finis Terrae II (FT2) uses Slurm for Resource Management.

- Detected issues when deploying Distributed TF on FT2:
  1. Distributed TF needs **Network Addresses (IP:Port)** in advance.

```
ServerDictionary={  
    "ps":["hostname01:2222","hostname02:2222"],  
    "worker":["hostname03:2222","hostname04:2222","hostname05:2222"]  
}
```



# Distributed TF: Queue system issues.

Distributed TF is based on communication protocol called gRPC.

CESGA Finis Terrae II (FT2) uses Slurm for Resource Management.

- Detected issues when deploying Distributed TF on FT2:
  1. Distributed TF needs **Network Addresses (IP:Port)** in advance.

```
ServerDictionary={  
    "ps":["hostname01:2222","hostname02:2222"],  
    "worker":["hostname03:2222","hostname04:2222","hostname05:2222"]  
}
```

**Queue system does not provide it.**



# Distributed TF: Queue system issues.

Distributed TF is based on communication protocol called gRPC.

CESGA Finis Terrae II (FT2) uses Slurm for Resource Management.

- Detected issues when deploying Distributed TF on FT2:
  1. Distributed TF needs **Network Addresses (IP:Port)** in advance.  
**Queue system does not provide it.**
  2. Distributed TF **Parameter Servers** run forever.



# Distributed TF: Queue system issues.

Distributed TF is based on communication protocol called gRPC.

CESGA Finis Terrae II (FT2) uses Slurm for Resource Management.

- Detected issues when deploying Distributed TF on FT2:
  1. Distributed TF needs **Network Addresses (IP:Port)** in advance.  
**Queue system does not provide it.**
  2. Distributed TF **Parameter Servers** run forever.  
**HPC Resources are wasted.**  
**User or Queue system have to stop them.**



# Solution: tf4slurm

CESGA has developed **tf4slurm** Python Package to solve the issues.

Python Module	Solved Issue
<b>ServerDictionary</b>	<b>IP:Port information</b>
<b>DistributedTFQueueHook</b>	<b>Close Distributed TF server gracefully*</b>

\* Adapted from <https://gist.github.com/yaroslavvb/82a5b5302449530ca5ff59df520c369e>

**Technical Report:**

<https://www.cesga.es/es/biblioteca/downloadAsset/id/803>

**GitHub repository:**

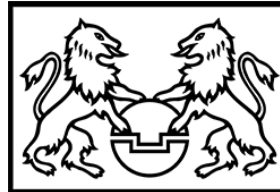
<https://github.com/gonfeco/tf4slurm>



# tf4slurm: test with an Industrial Case

## Fortissimo H2020 Project

### Experiment 707: Cyber-Physical Laser Metal Deposition (CyPLAM)

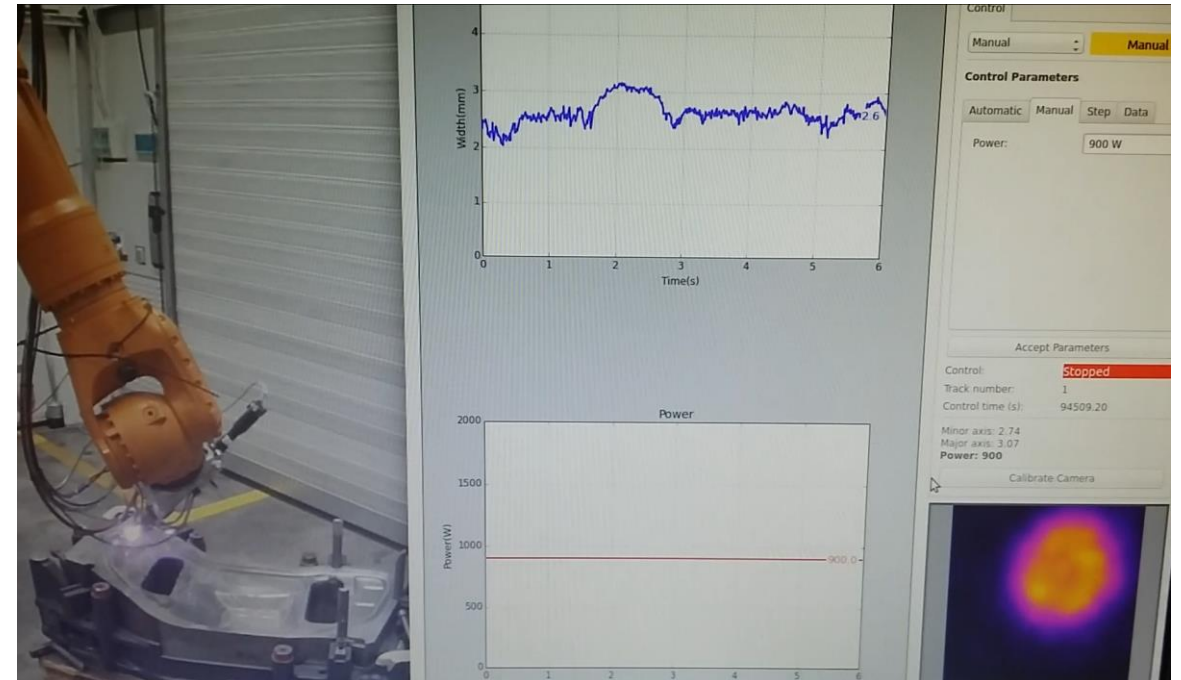


# CyPLAM description

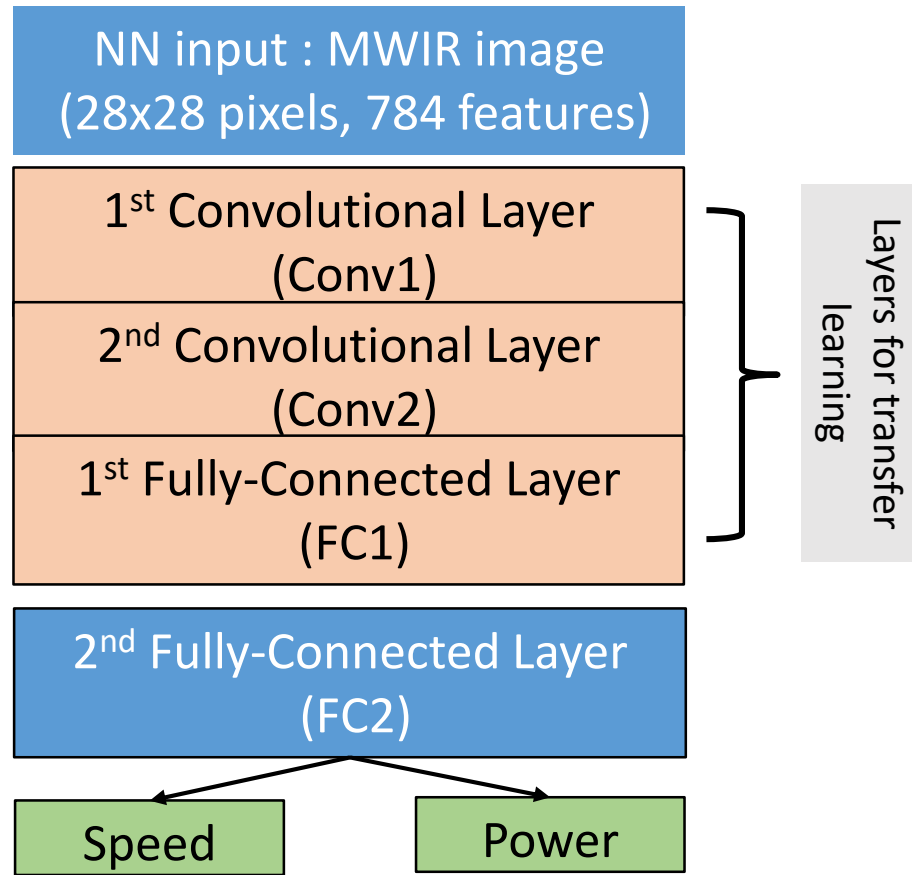
Using Laser Metal Deposition (LMD) for building and repairing large metal parts.

LMD process recorder by Medium Wavelength Infrared (MWIR) sensors attached to laser header.

Use ML algorithms for monitoring the LMD process based on the MWIR images.

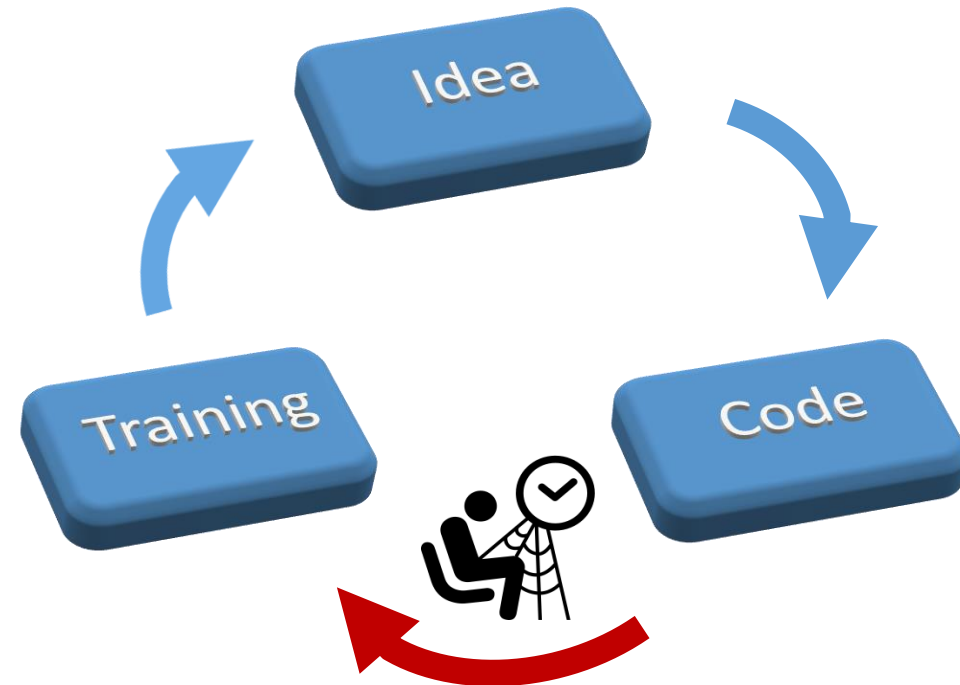


# CyPLAM Algorithm



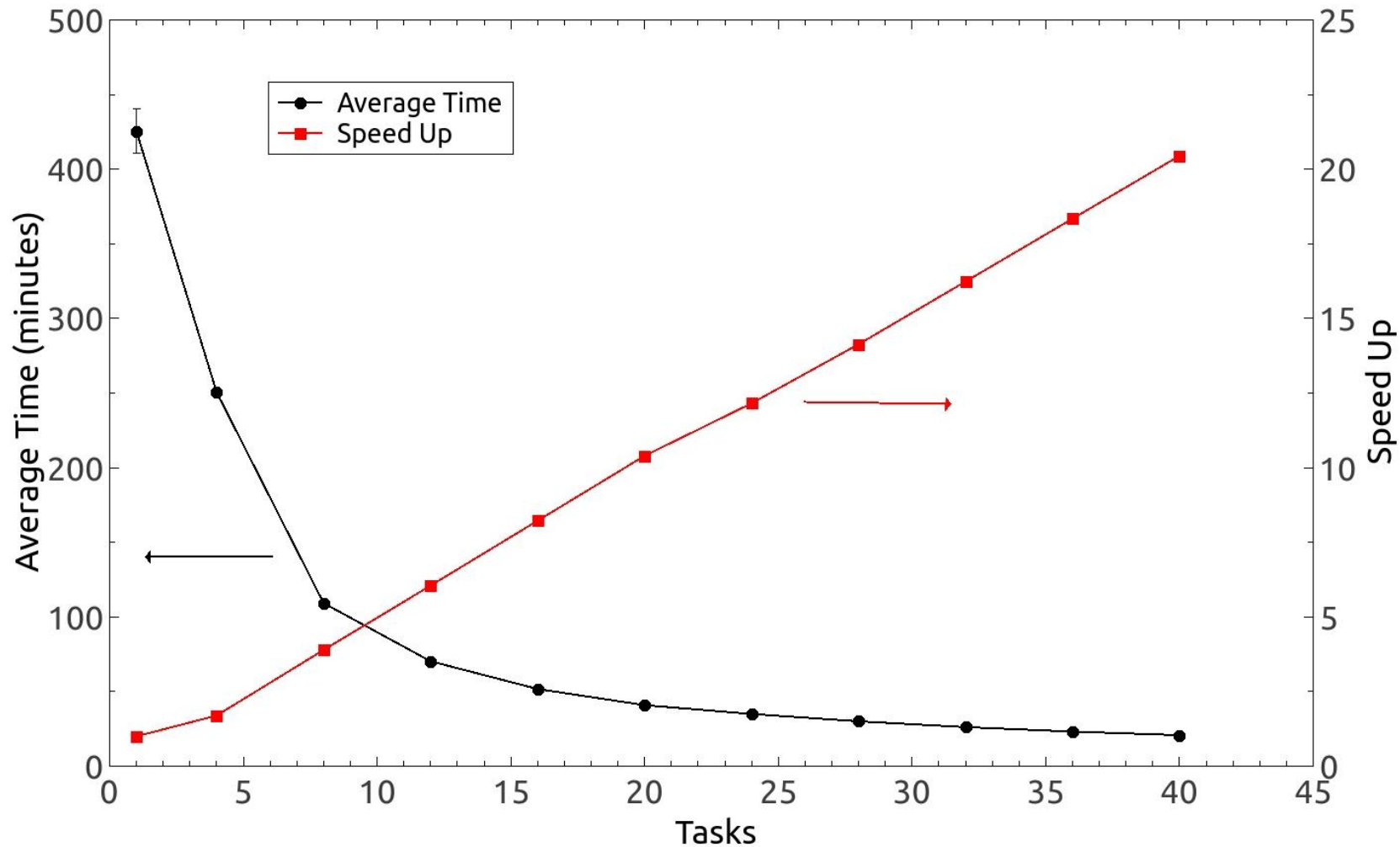
NN Graph model (extracted from Tensorboard)

**For Biggest tested model:  
Training Time > 7 h.**

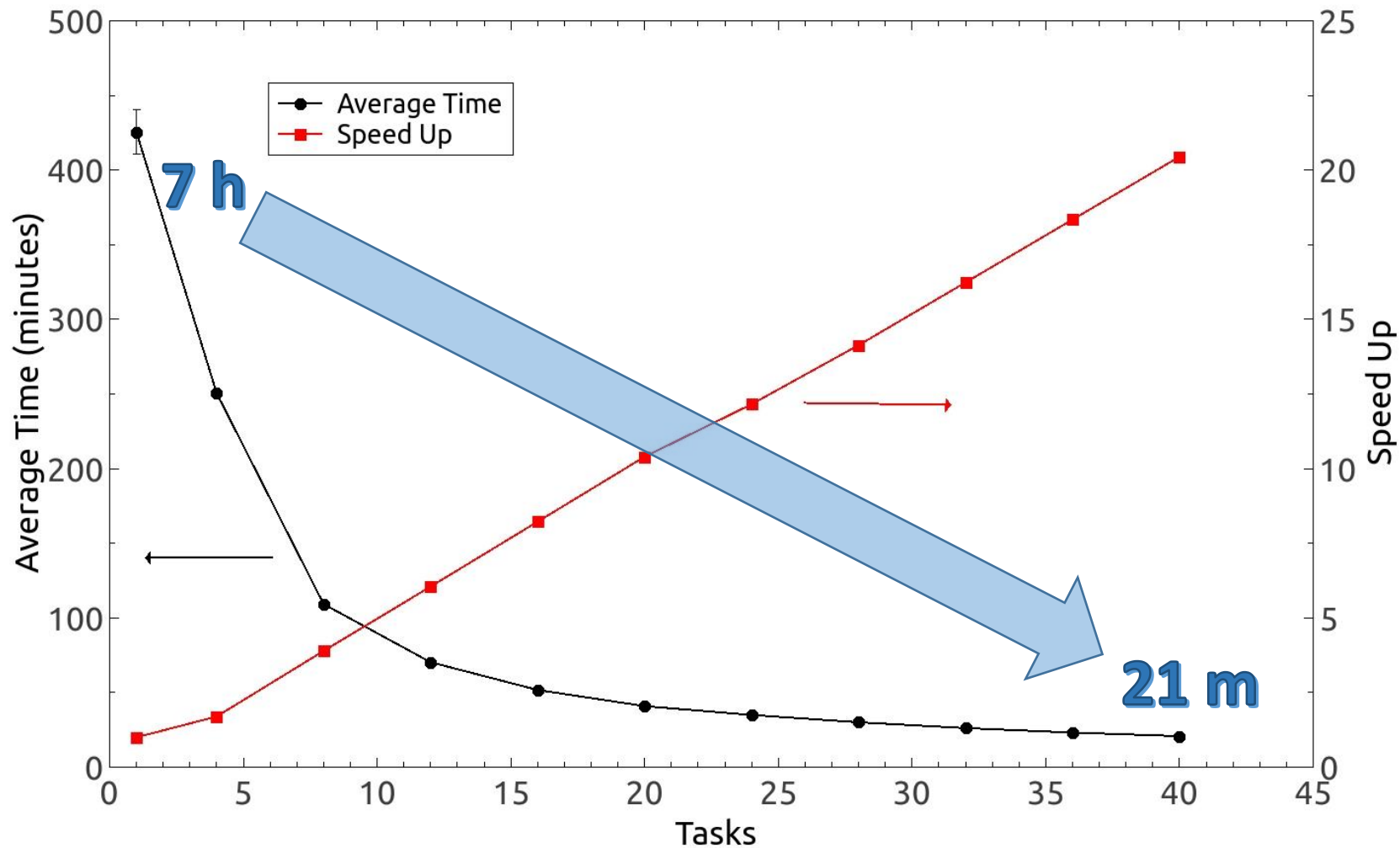




# CyPLAM training using tf4slurm



# CyPLAM training using tf4slurm



**7 hours**  
↓  
**~20 minutes**



# Summary and Conclusions.

- Several issues deploying Distributed TensorFlow on CESGA Finis Terrae II were detected:
  - Queue system does not provide **mandatory IP:Port** information in advance.
  - HPC Resources are **wasted** due to **Parameter Servers running forever**.
- **tf4slurm** Python Package was developed to **solve** these detected issues.
- **tf4slurm** was tested using an Industrial Case:
  - Largest training **reduced from 7 hours to near 20 minutes**.
- HPC can **greatly decrease the design time** of ML algorithms **boosting productivity**.





FORTISSIMO



# THANKS FOR YOUR ATTENTION !!!

## Technical Report:

- <https://www.cesga.es/es/biblioteca/downloadAsset/id/803>

## GitHub repository:

- <https://github.com/gonfeco/tf4slurm>

