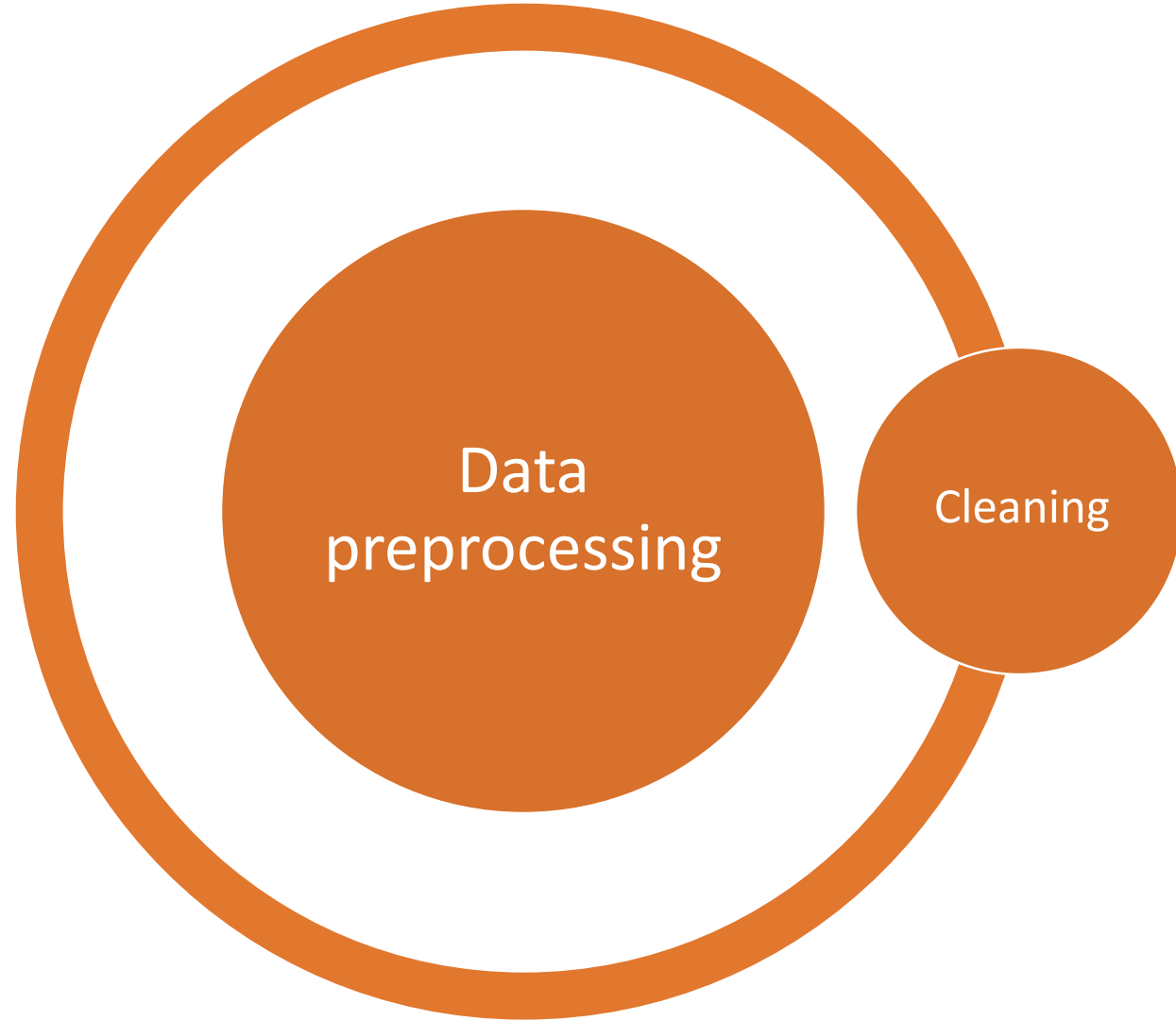




Data Preprocessing for DL

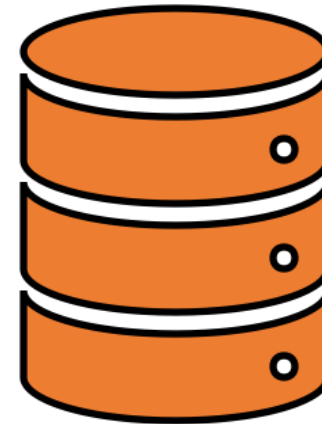
Ing. Ján Pavlus

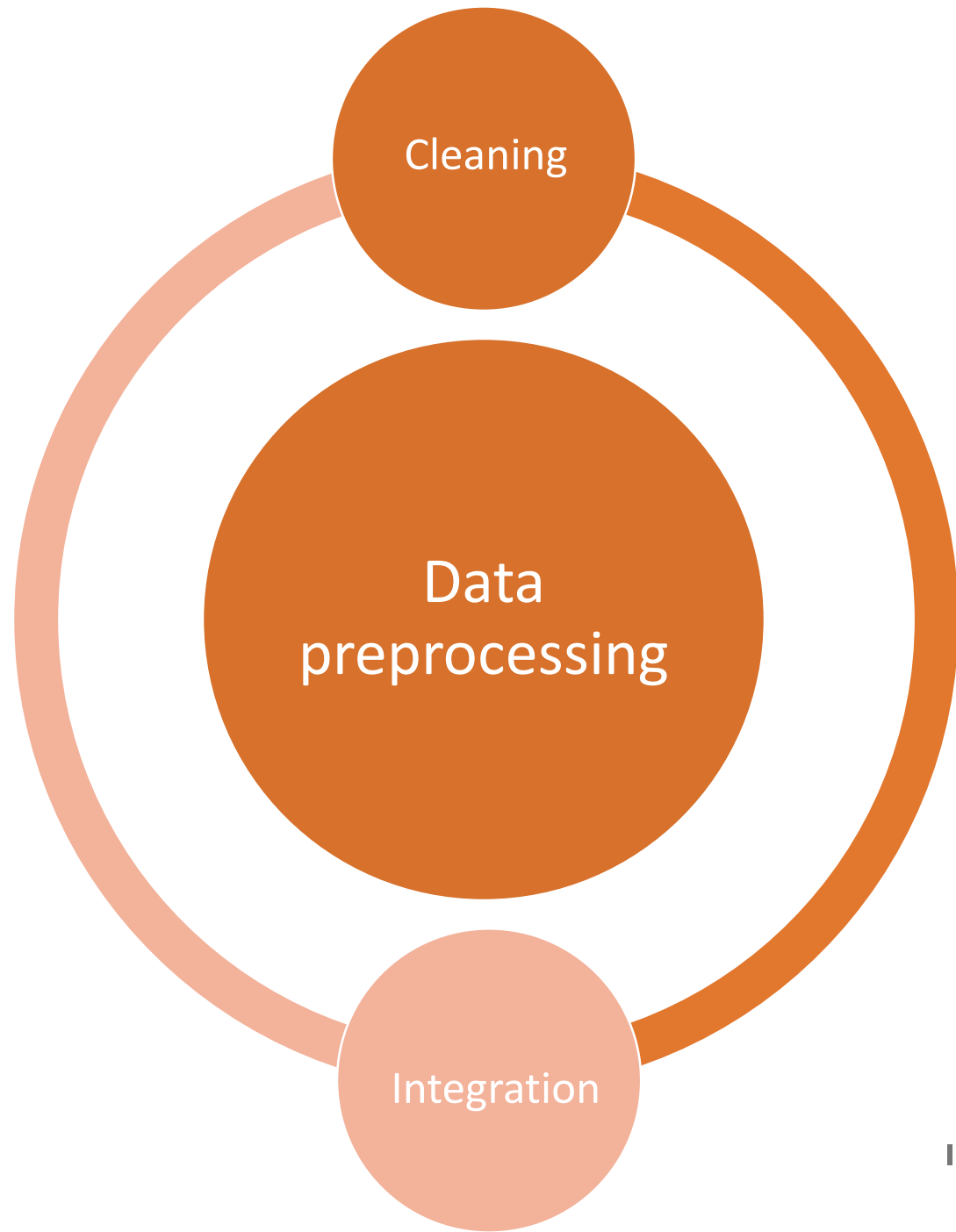


Cleaning

Data cleaning

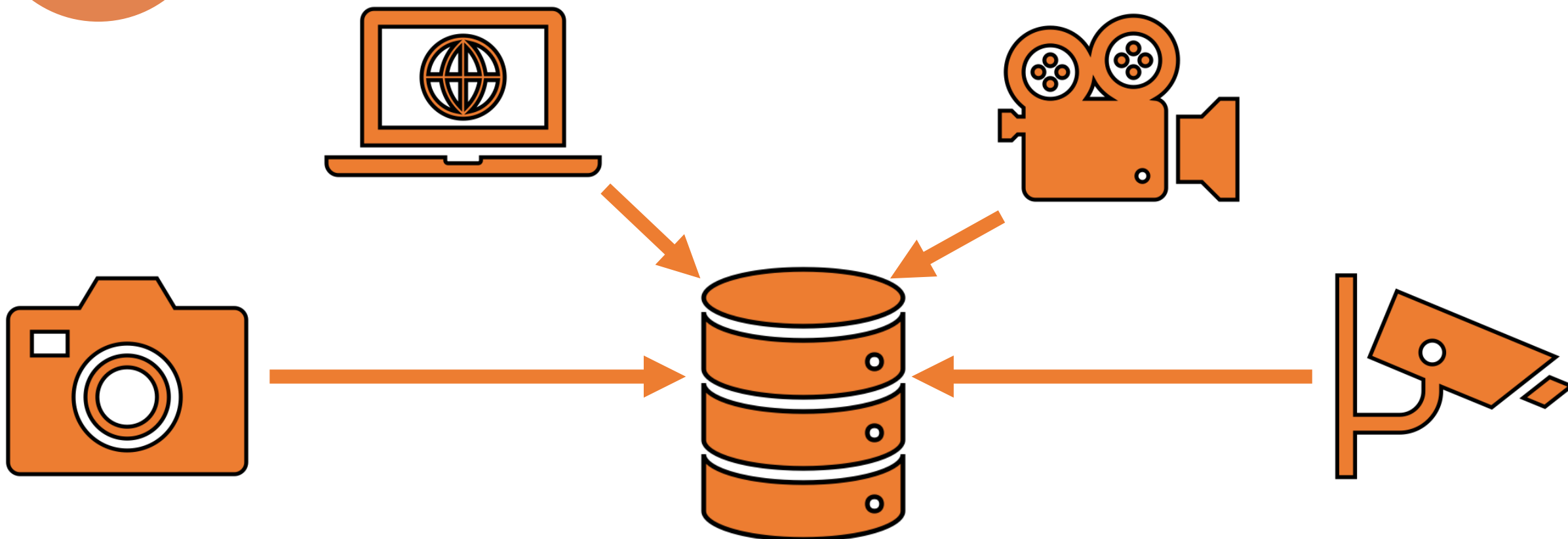
- Missing values
- Noisy data
- Removing outliers
- Duplicities

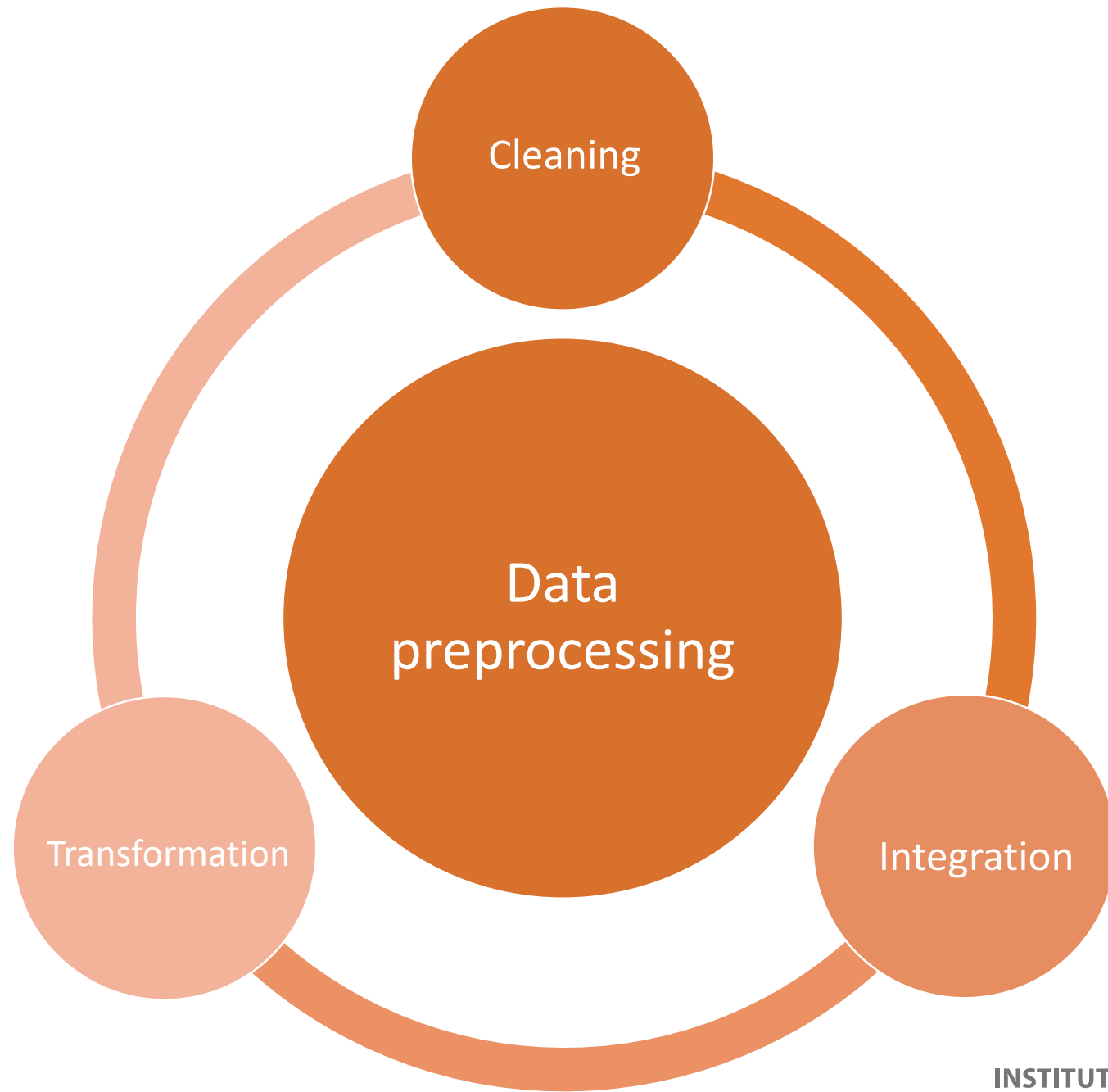




Integration

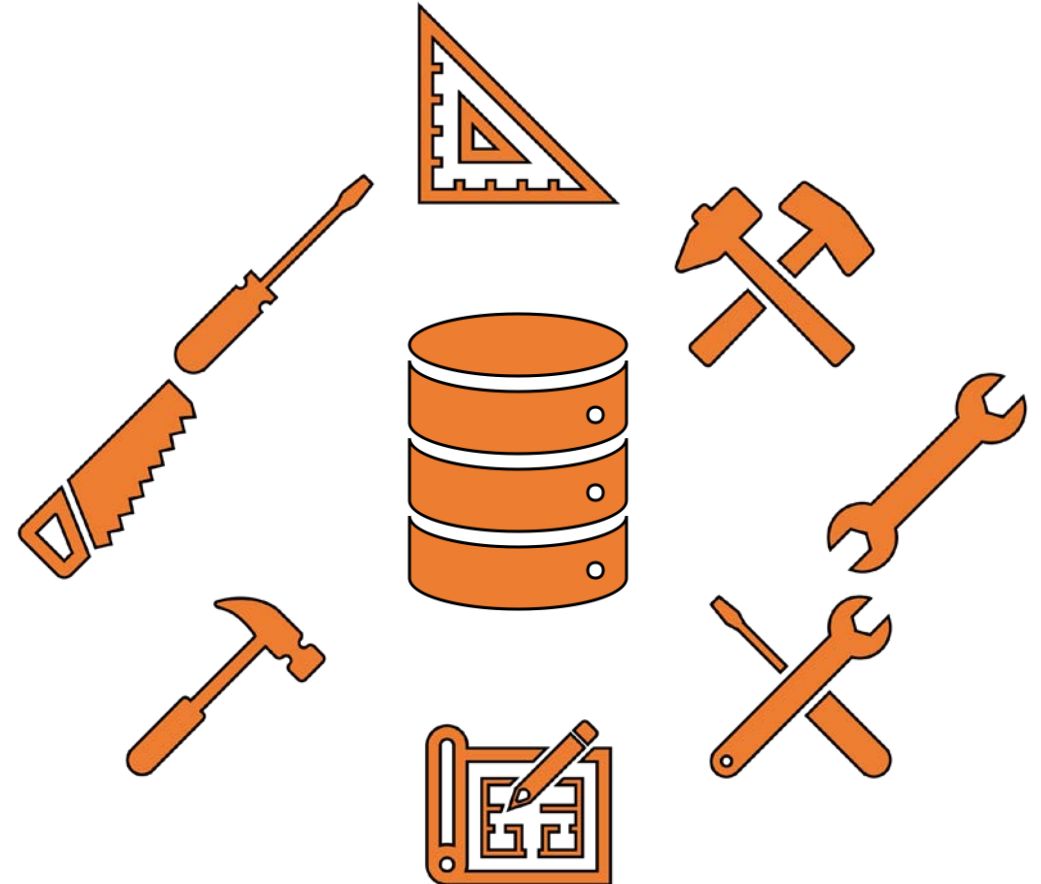
Data integration





Data transformation

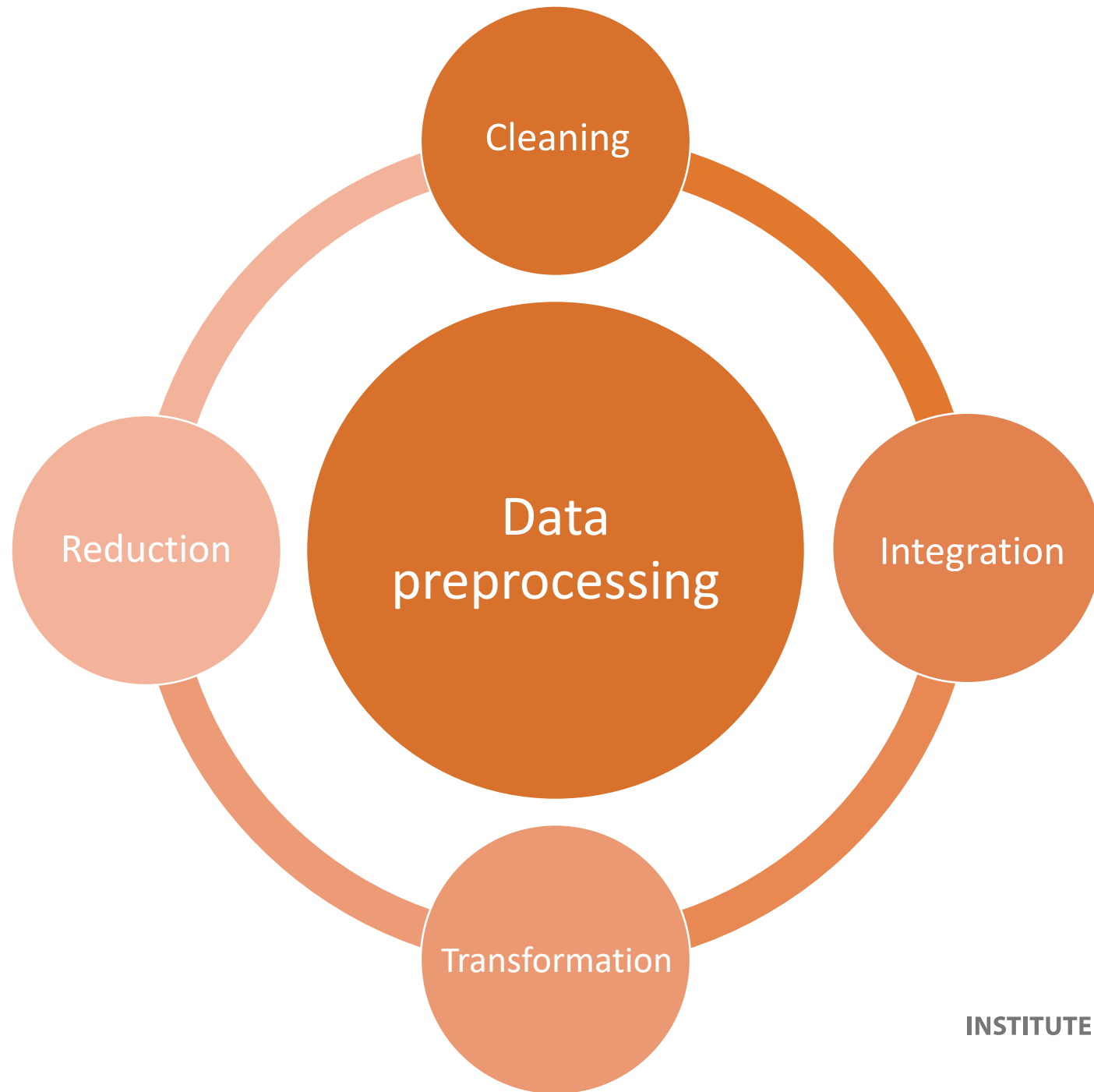
- Normalization
- Attribute selection
- Augmentations



Data transformation - Augmentations

- Extend dataset
- Extend variability
- Robust system
- Prevent from overfitting

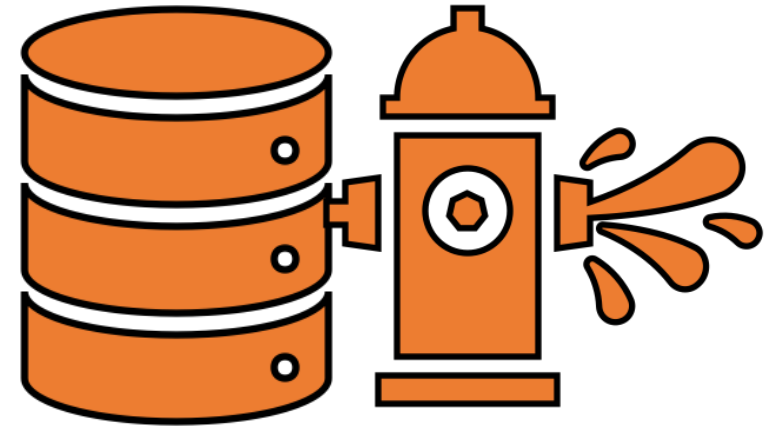


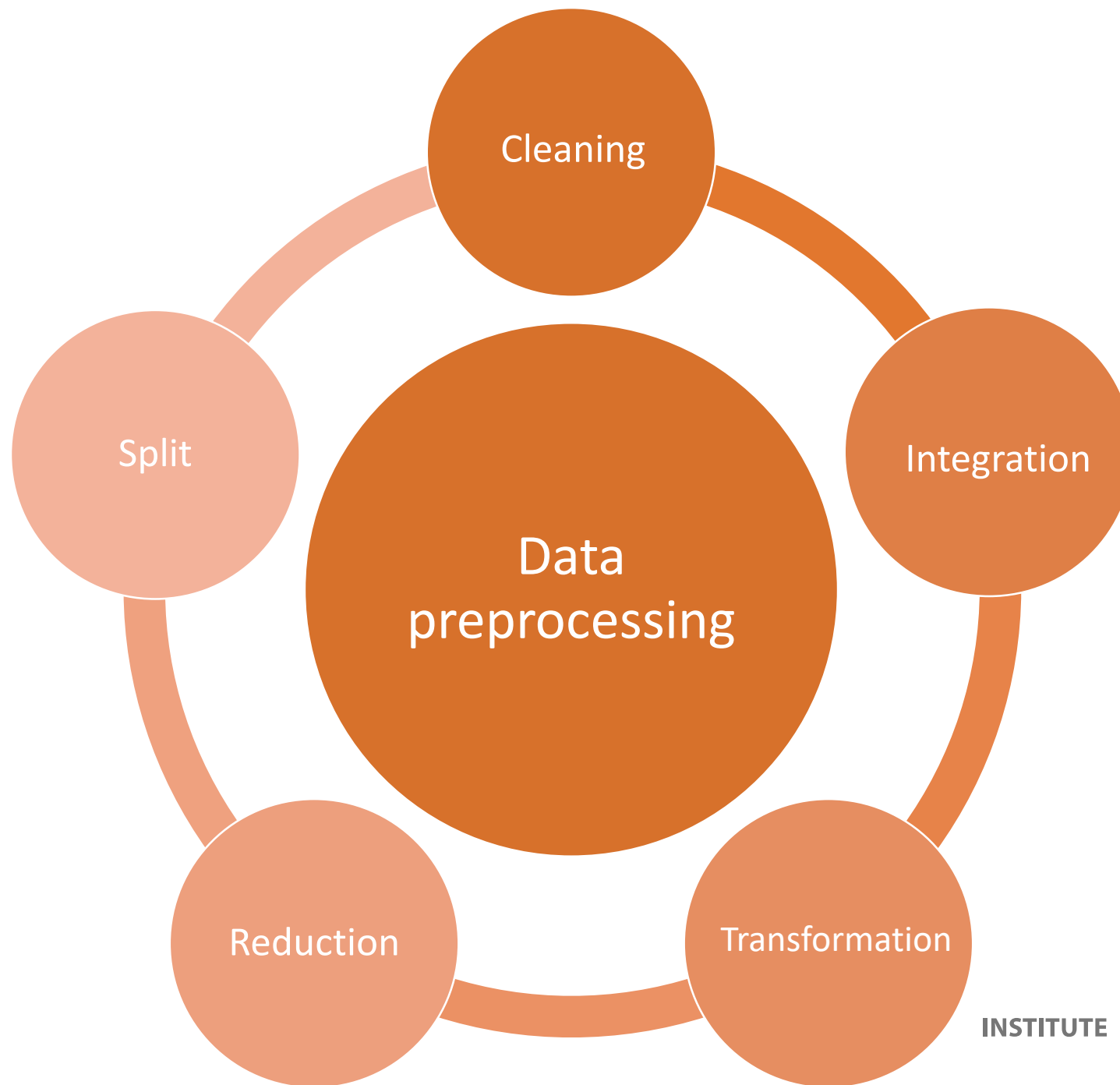


Reduction

Data reduction

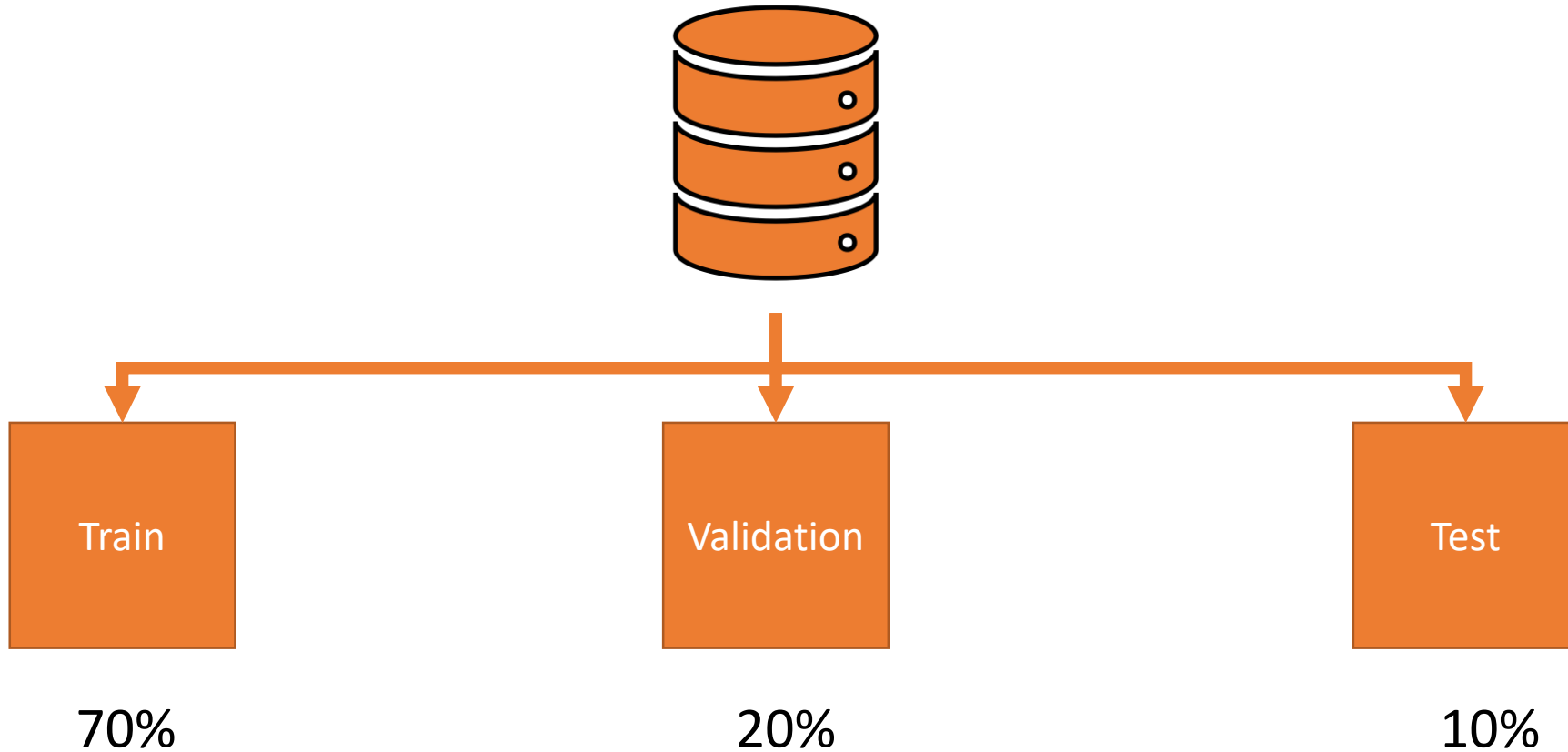
- Compress data
- Choose right representation
- Remove unnecessary attributes

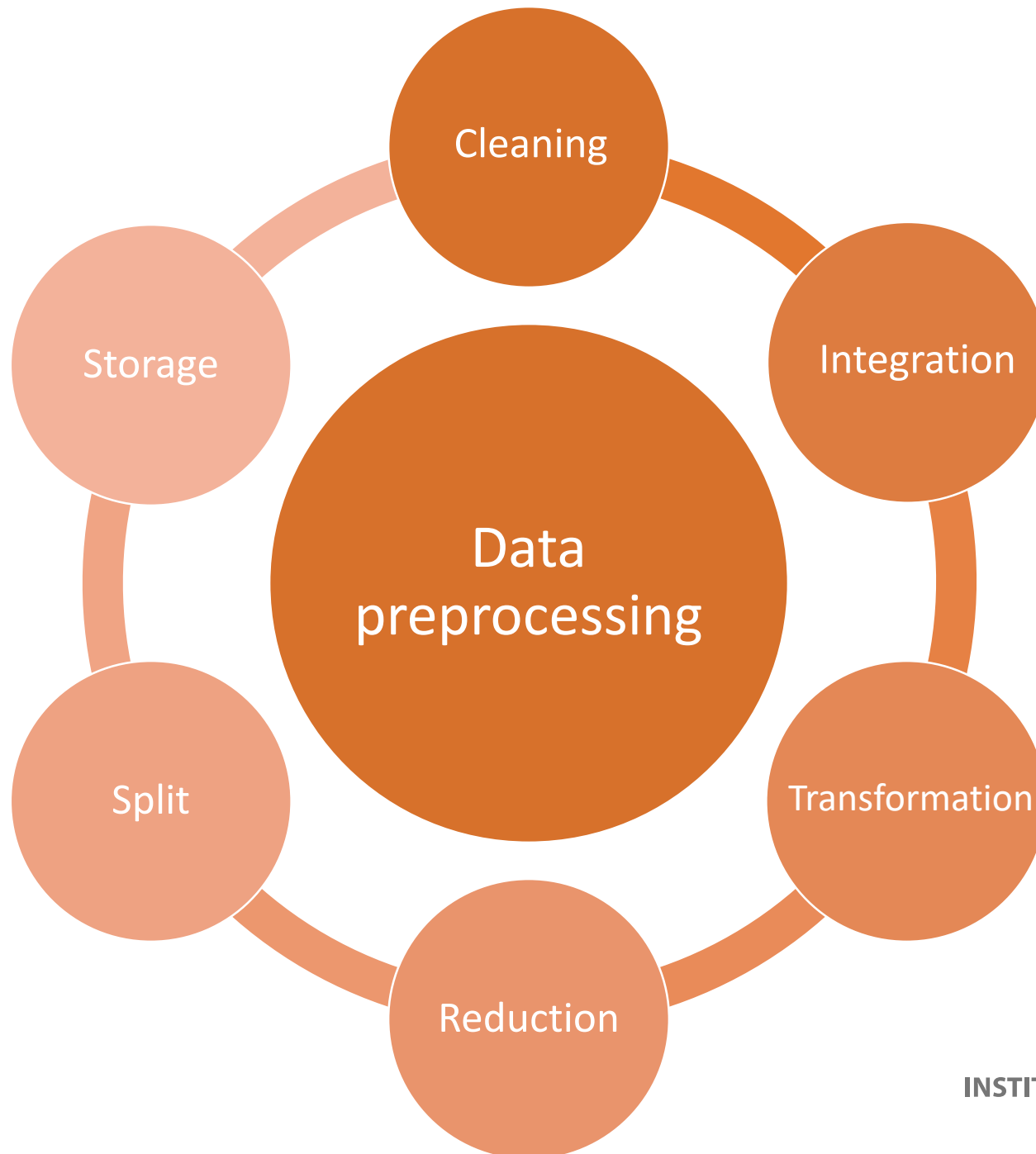




Storage

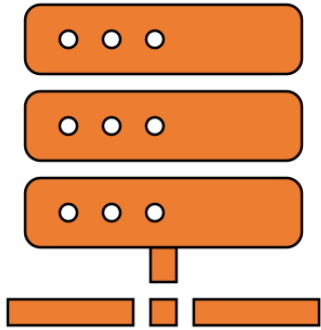
Data split



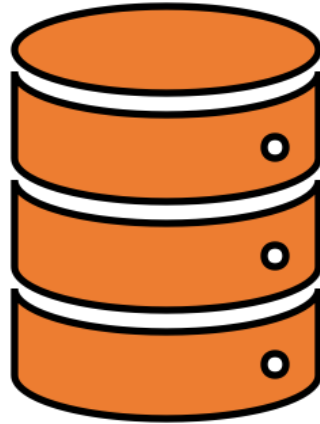


Storage

Data storage

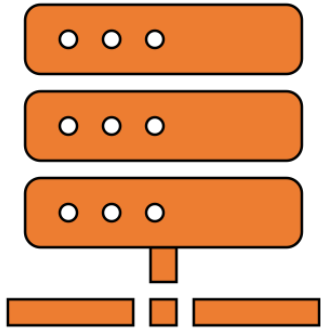


HDD



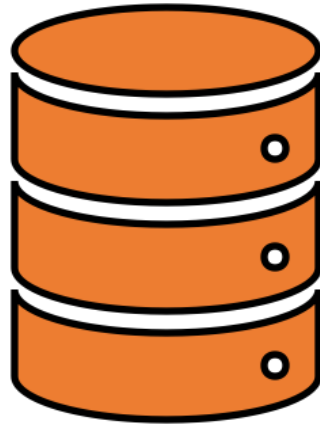
Storage

Data storage



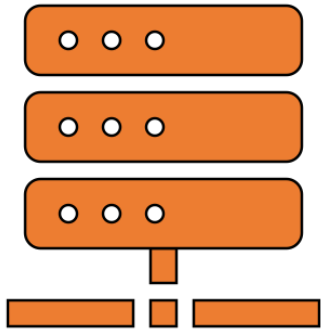
HDD

+ Cheap



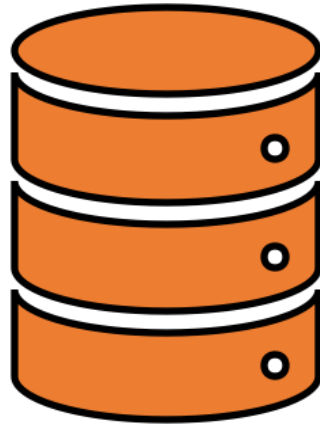
Storage

Data storage



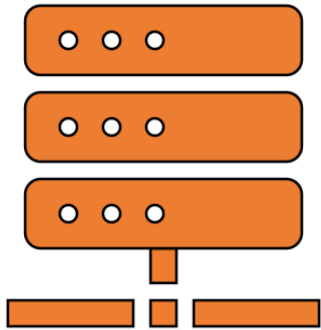
HDD

- + Cheap
- + Store big data



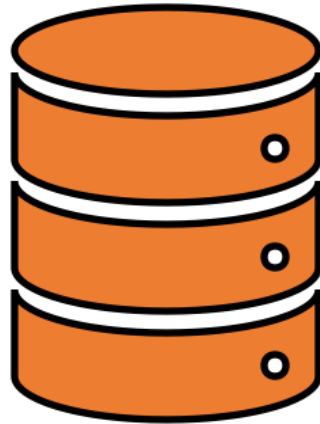
Storage

Data storage



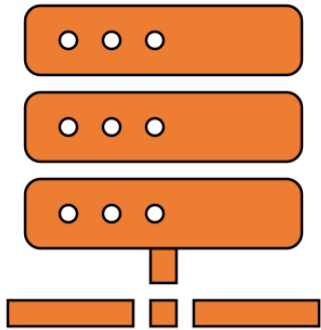
HDD

- + Cheap
- + Store big data
- Slow



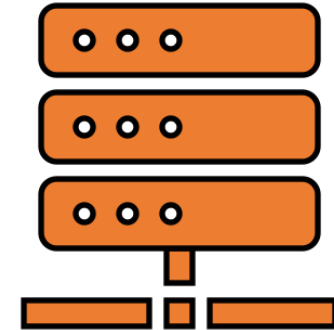
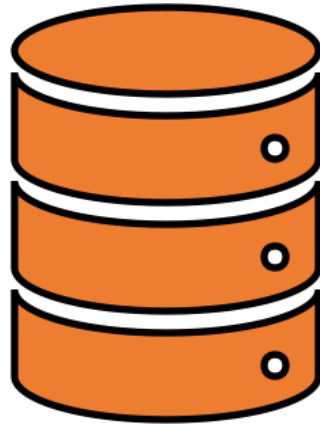
Storage

Data storage



HDD

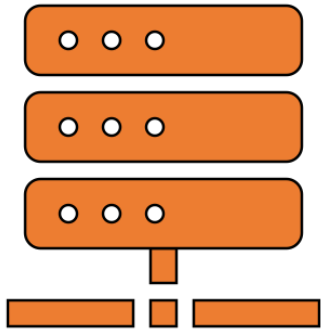
- + Cheap
- + Store big data
- Slow



SSD

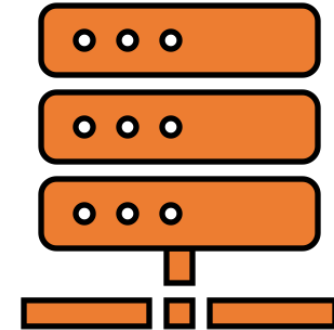
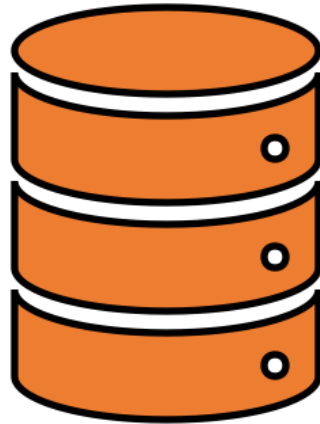
Storage

Data storage



HDD

- + Cheap
- + Store big data
- Slow

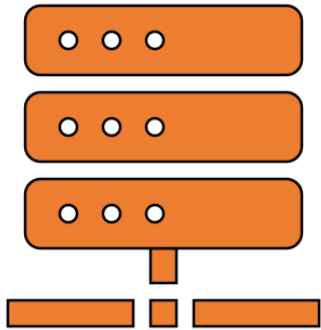


SSD

- + Really fast

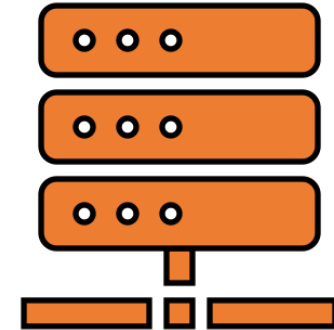
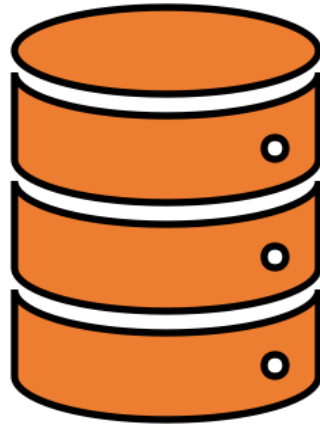
Storage

Data storage



HDD

- + Cheap
- + Store big data
- Slow

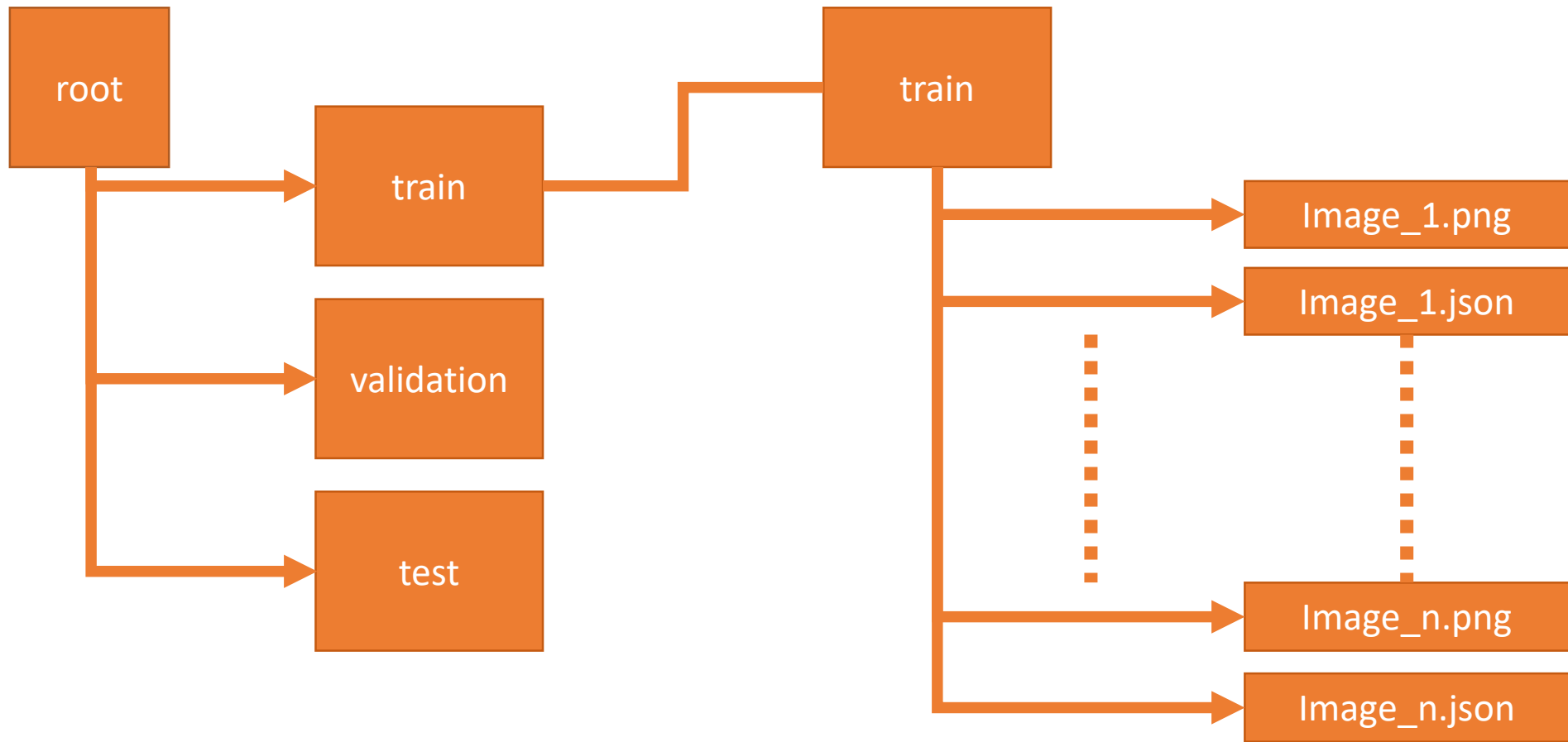


SSD

- + Really fast
- Expensive

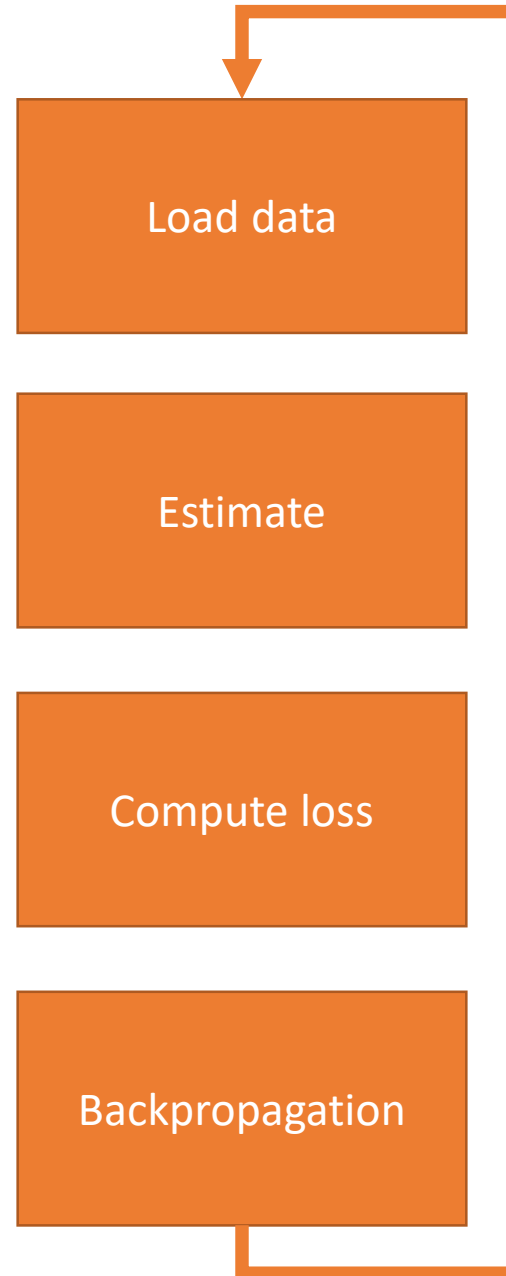
Storage

Data storage – folder structure



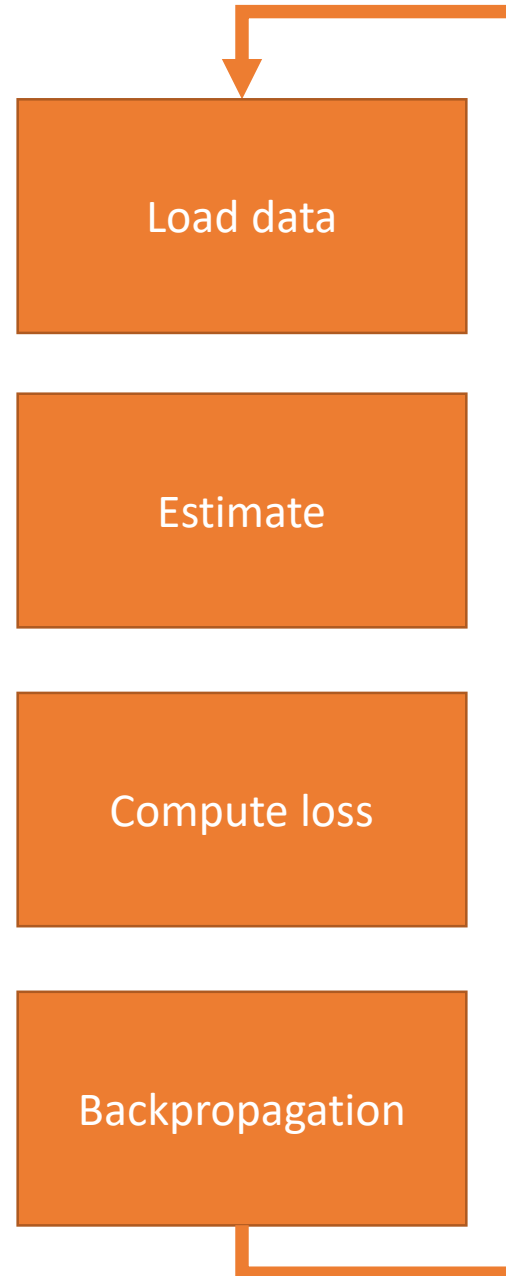
Dataloader

- Load data



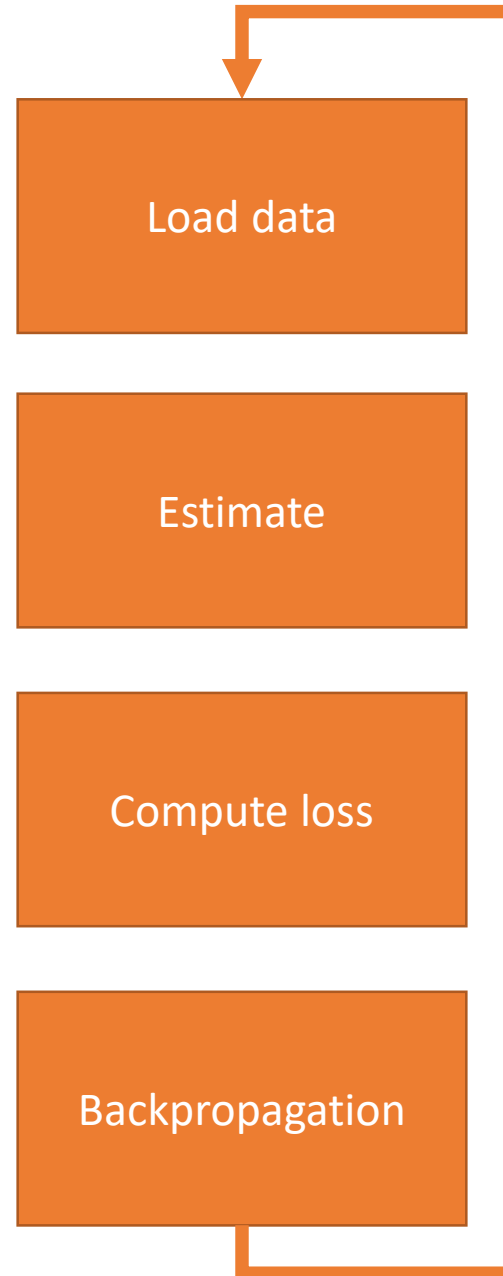
Dataloader

- Load data
- Multithread



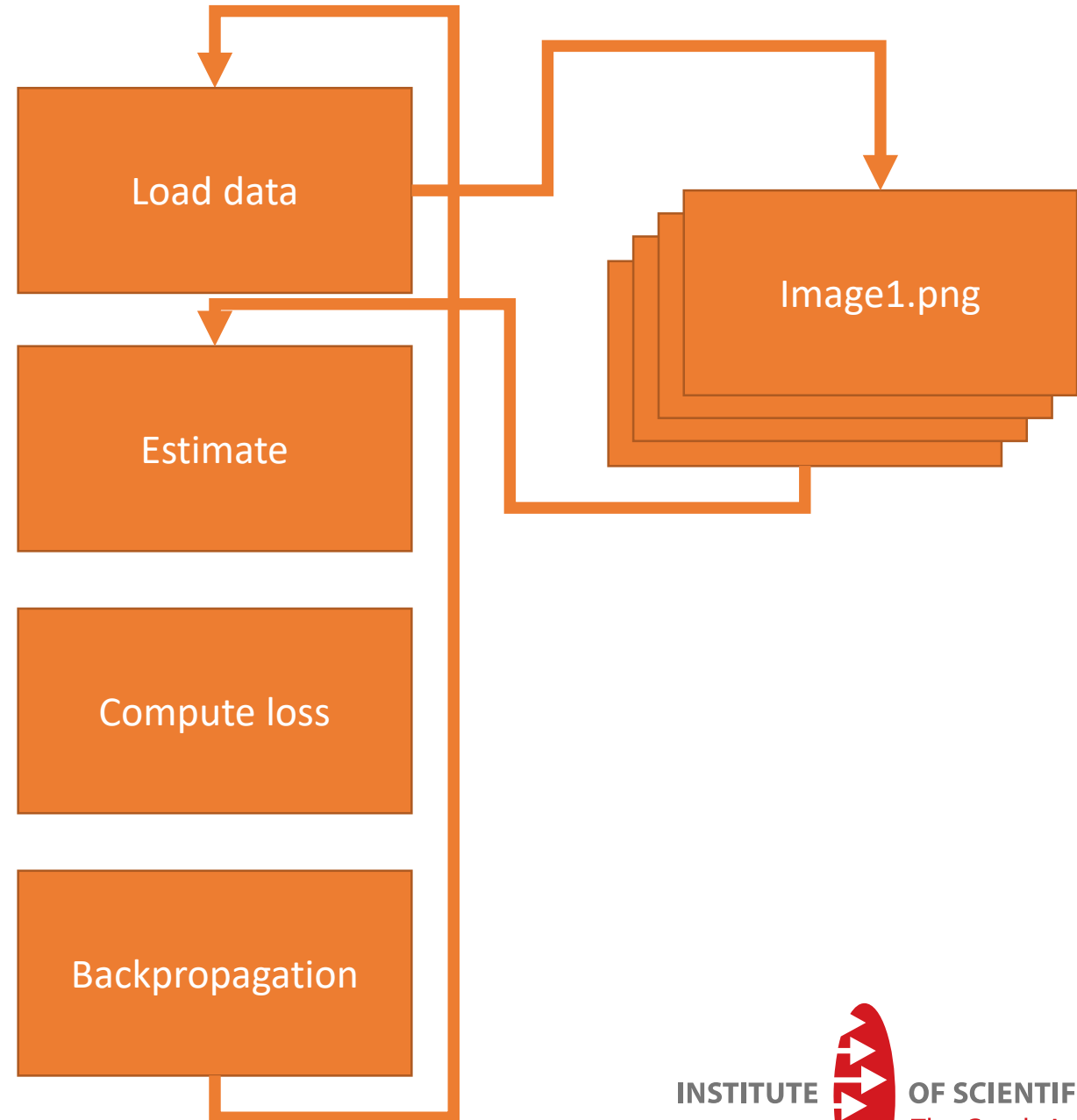
Dataloader

- Load data
- Multithread
- Provides data for the neural network



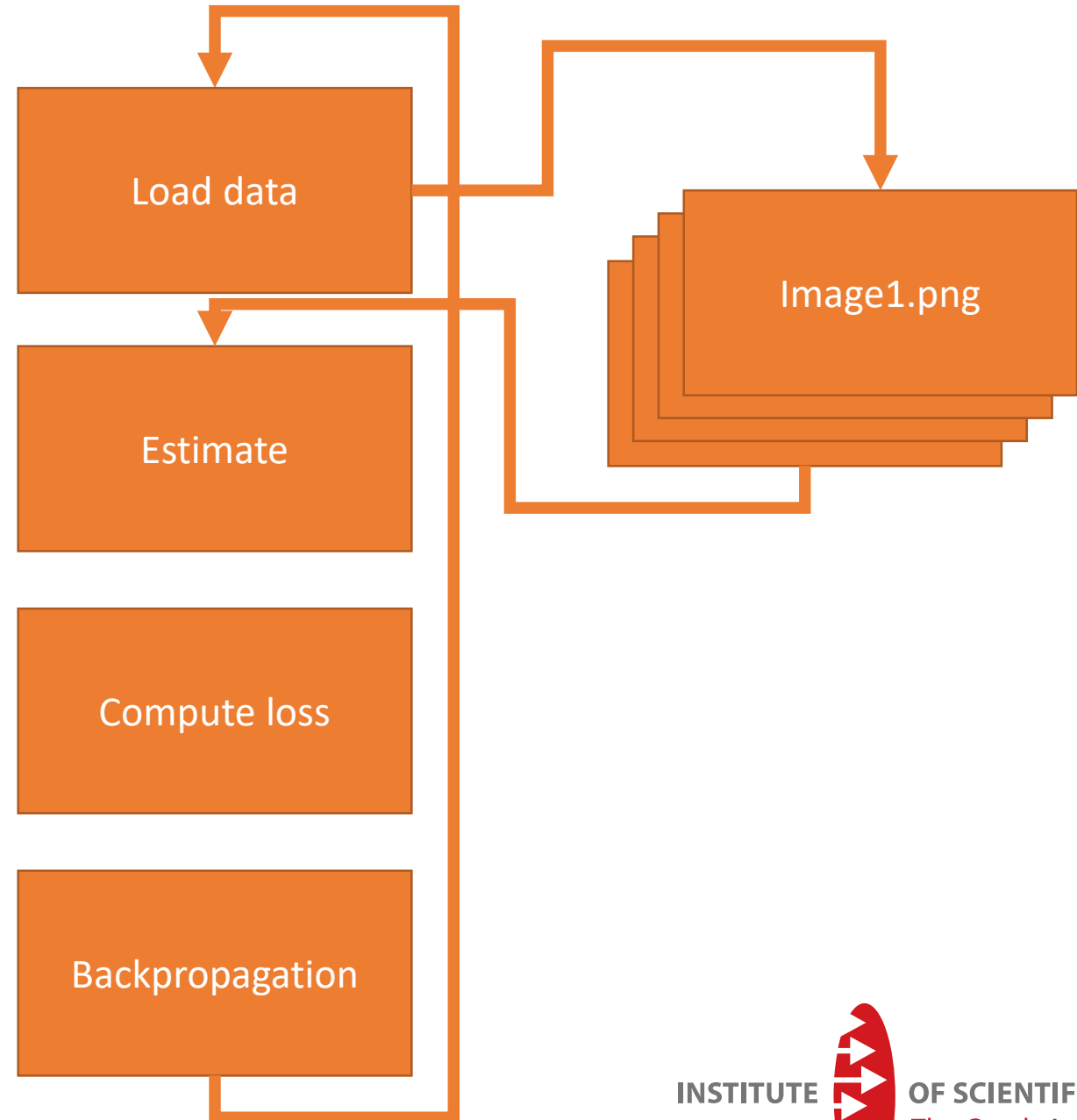
Dataloader

- Load data
- Multithread
- Provides data for the neural network
- Create batches



Dataloader

- Load data
- Multithread
- Provides data for the neural network
- Create batches
- Chunk data



Dataloader

- Load data
- Multithread
- Provides data for the neural network
- Create batches
- Chunk data
- Make online augmentation

