

aidge

The first independent, open-source platform
dedicated to embedded AI



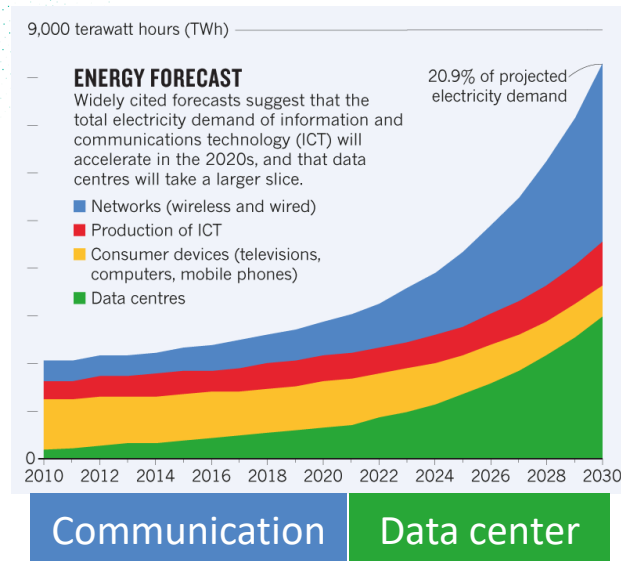
Embedded AI : deep digital transformation

OPPORTUNITY

- Real time
- Data and Model Security
- Cost reduction

NECESSITY

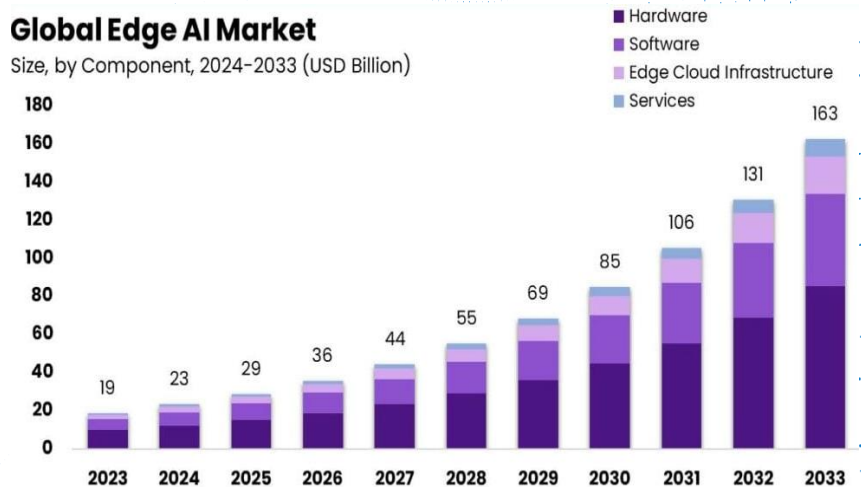
- Scaling up AI
- Economic Challenge
- Strategic challenge



[Nature]

Global Edge AI Market

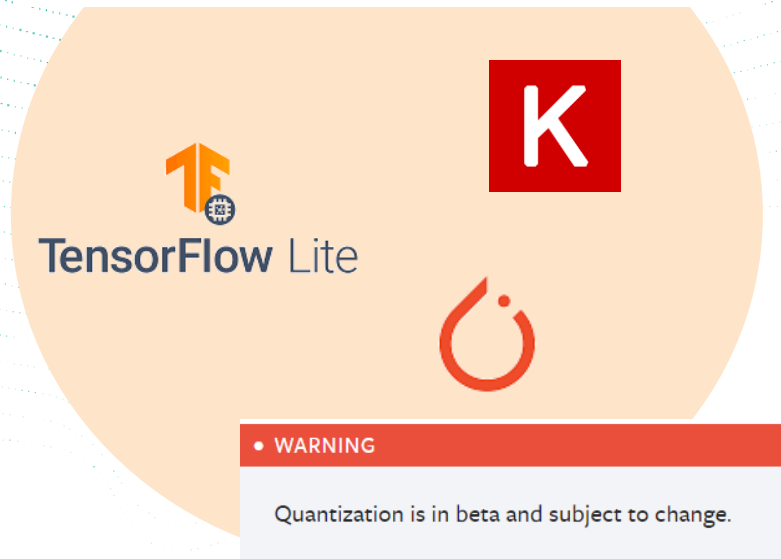
Size, by Component, 2024-2033 (USD Billion)



[Market.us]

Existing tools : strong orientations

Deep Learning Platforms



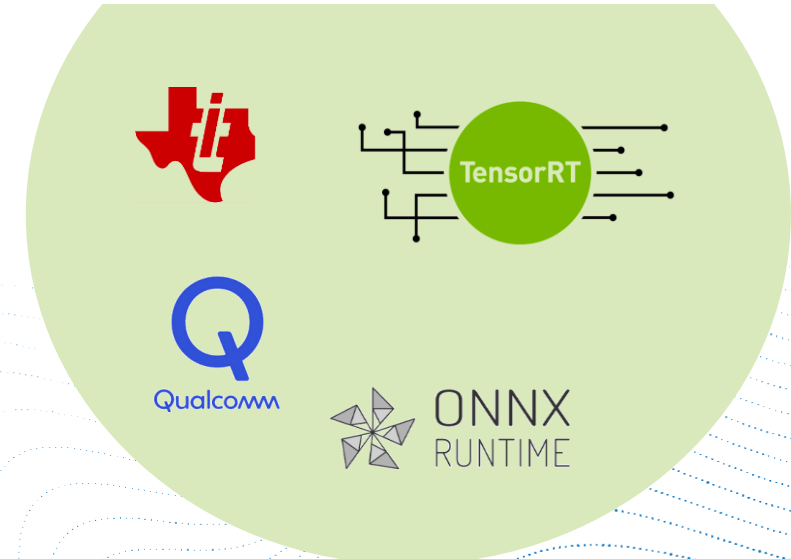
Close source / Maturity
Dependencies

Compilers



Low level optimization
Black box

Hardware SDK



Close source
Hardware specific

Challenges

OPENESS

Reusable and
adaptable tools to
foster innovation

INDEPENDENCE

Free choice of
material from
component to
system

INNOVATION

Integrating innovative
paradigms for trust
and frugality

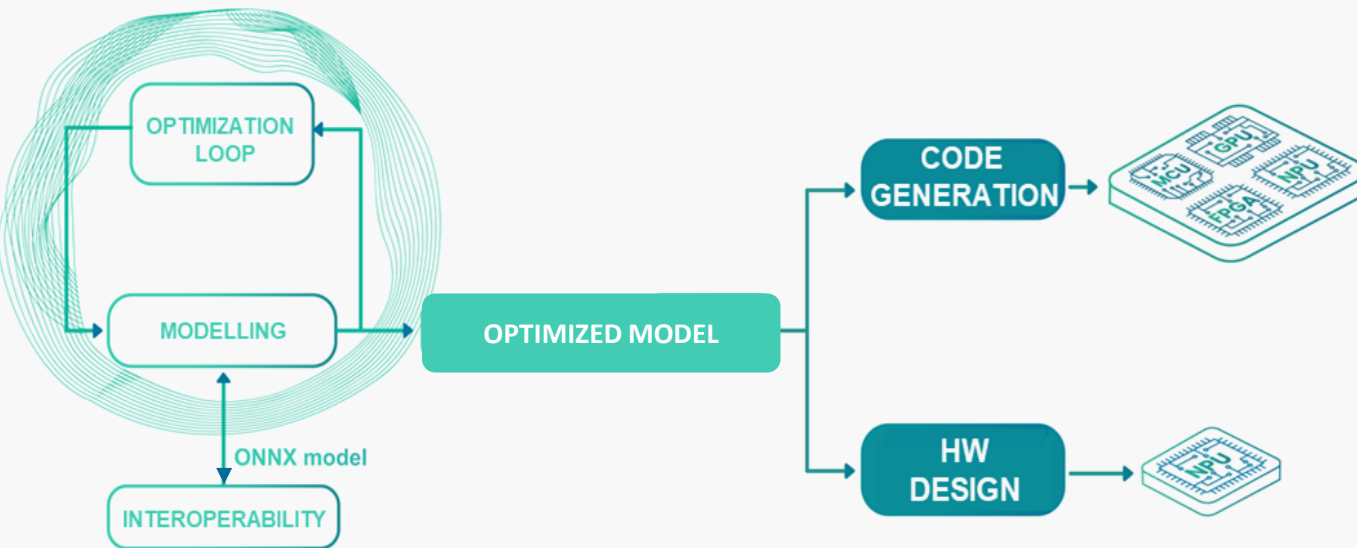
Complementary consortium

aidge



aidge

The first open, independent platform dedicated to embedded AI



ONNX

PyTorch TensorFlow K

Hosted by **ECLIPSE**
FOUNDATION



Integrated platform

- From import to deployment
- High degree of interoperability
- Minimal dependencies



Modular and expandable platform

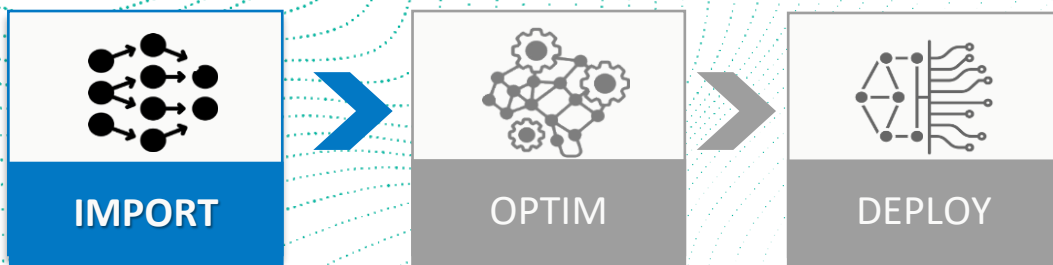
- Lightweight core module with plug-ins
- Open-source collaborative environment

API



Multi-Platform and Packaging





- **High degree of interoperability** with ONNX standard

+60 operators and involved in the Safety ONNX standard

- **Native support of the main embedded architectures**

CNN, RNN, GAN, YOLO, Transformer and soon SNN

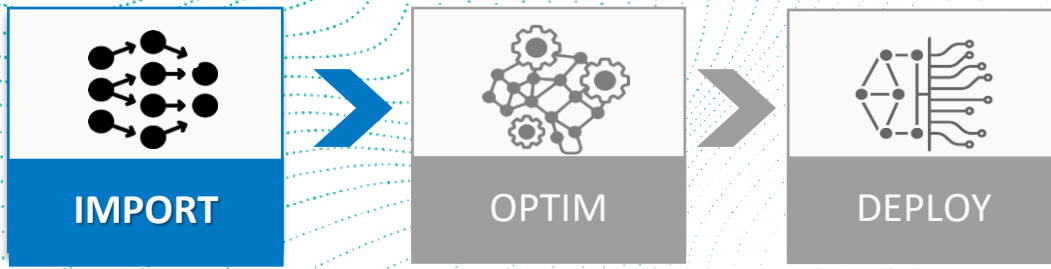
- Rich analysis tools to assess model complexity: parameters, operations, etc.
- Unique intermediate representation for easy model access and manipulation

**ONNX coverage ratio
DINOv2 (Meta): 100%.**

```
# Here show nb of operators!
aidge_onnx.native_coverage_report(dinov2_model)
```

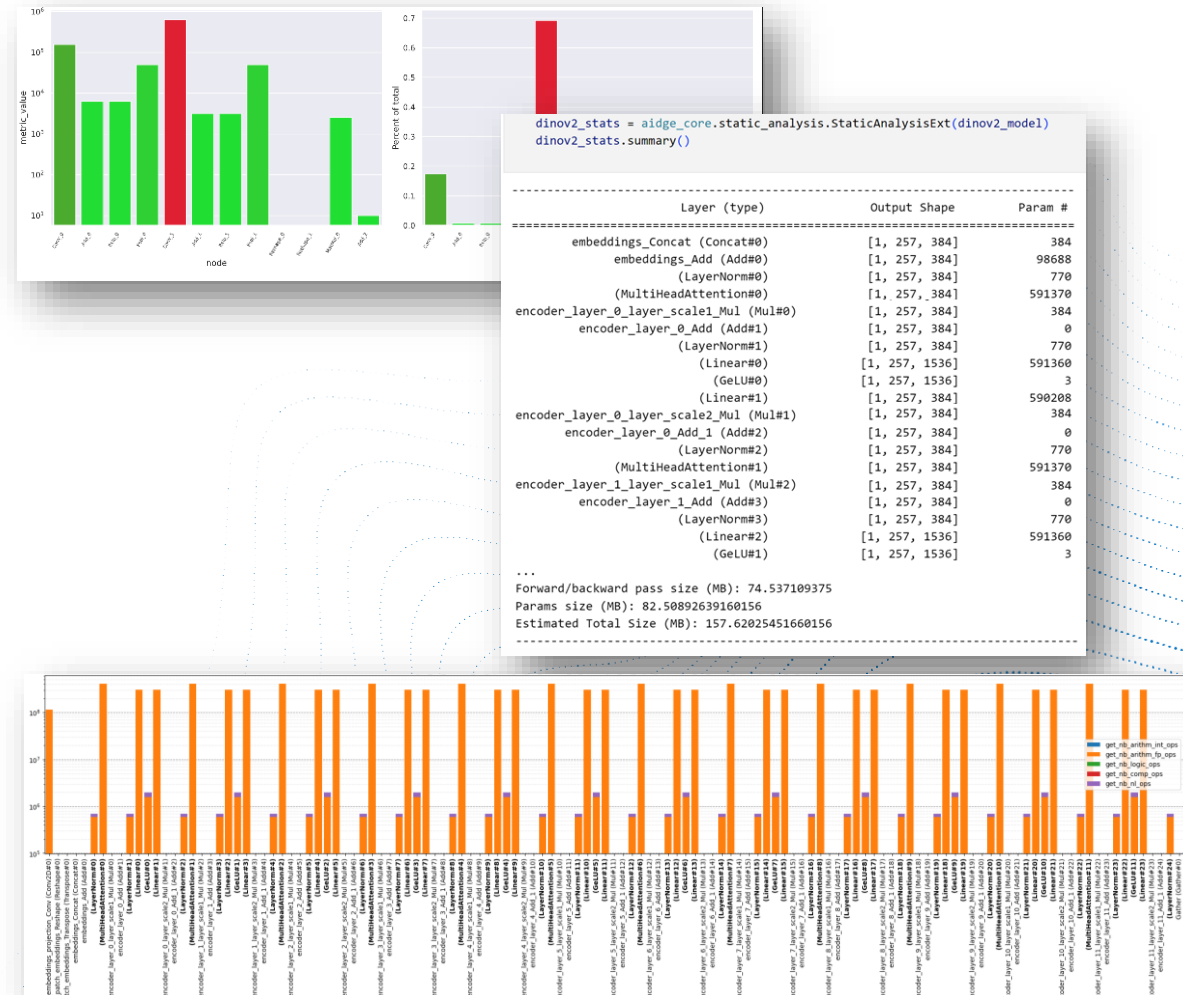
[4]

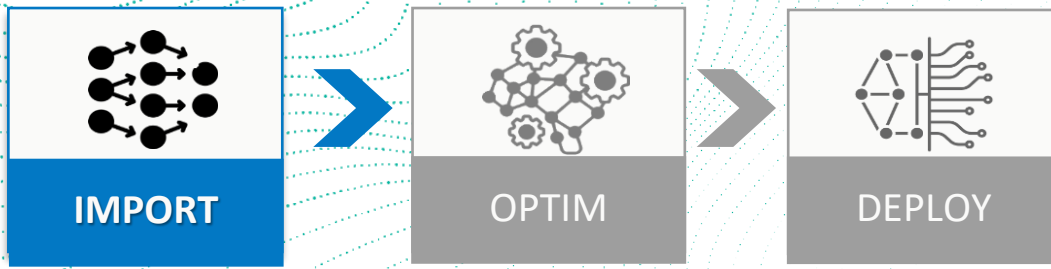
```
... Native operators: 824 (17 types)
- Add: 159
- Concat: 1
- Conv2D: 1
- Div: 49
- Erf: 12
- Gather: 1
- MatMul: 72
- Mul: 73
- Pow: 25
- Producer: 209
- ReduceMean: 50
- Reshape: 49
- Softmax: 12
- Split: 12
- Sqrt: 25
- Sub: 25
- Transpose: 49
Generic operators: 0 (0 types)
Native types coverage: 100.0% (17/17)
Native operators coverage: 100.0% (824/824)
```



aidge

- High degree of interoperability with ONNX standard
- Native support of the main embedded architectures
- **Rich analysis tools to assess model complexity:** parameters, operations, etc.
- Unique intermediate representation for easy model access and manipulation

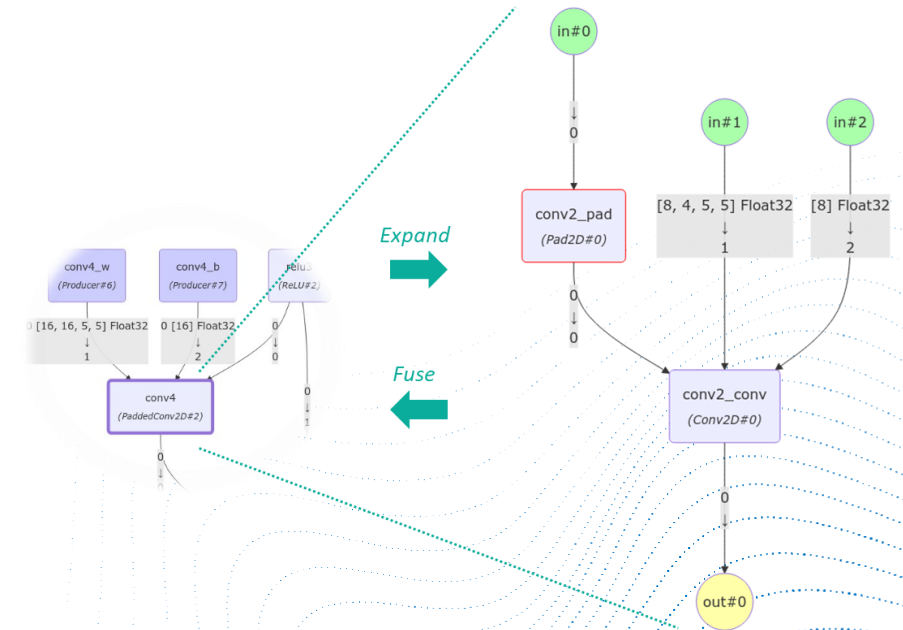




aidge

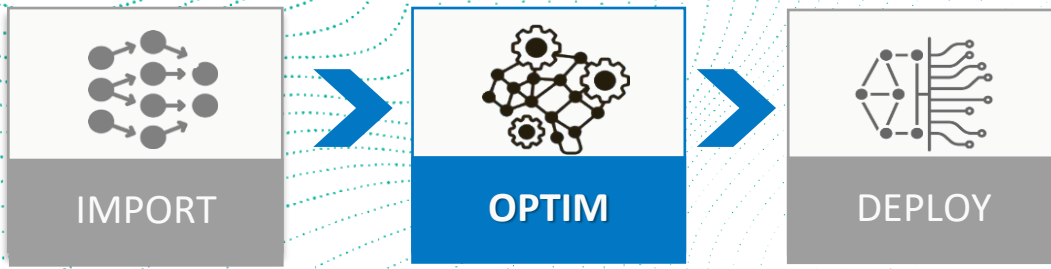
Match the granularity **required by the implementation**

- High degree of interoperability with ONNX standard
- Native support of the main embedded architectures
- Rich analysis tools to assess model complexity: parameters, operations, etc.
- **Unique intermediate representation** for easy model access and manipulation



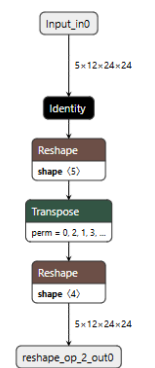
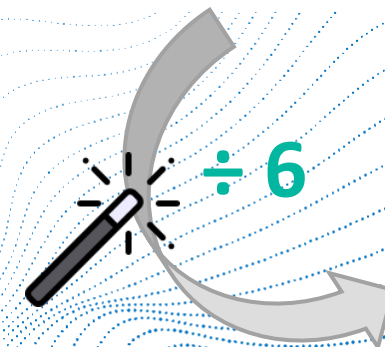
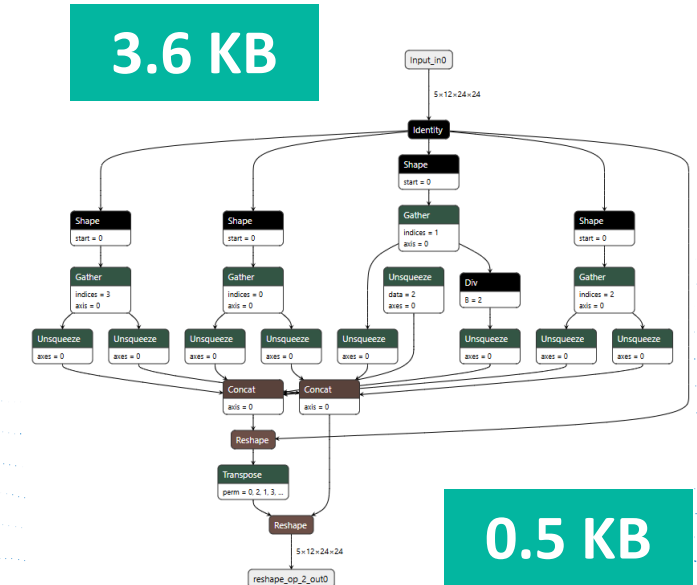
A powerful **graph matching system**

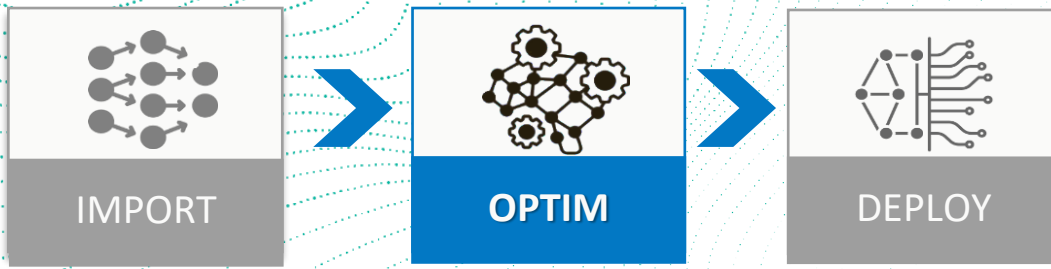
```
aidge_core.fuse_to_metaops(dinov2_model,
    "ScaledDotProductAttention#1->Transpose->Reshape#1->Linear;"
    "Reshape#1<1~Producer;"
    "ScaledDotProductAttention#1<0-(Transpose<-Reshape#2<-Add#1);"
    "ScaledDotProductAttention#1<1-(Transpose<-Reshape#3<-Add#2);"
    "ScaledDotProductAttention#1<2-(Transpose<-Reshape#4<-Add#3);"
    "Reshape#2<1~Producer;" "MultiHeadAttention")
```



aidge

- **Automatic model reduction:** catalog of optimizations with deletion, reorganization and merging of operations
- **State-of-the-art quantification to desired accuracy (ResNet, VGG, etc.)**
 - After learning: without loss up to 8-bit integer
 - During learning: without loss up to 4-bit integer
- **Tensor decomposition compression method**
 - ResNet-50 x ImageNet: 15% compression without loss
 - ResNet-18 x CIFAR 100: 45% compression without loss

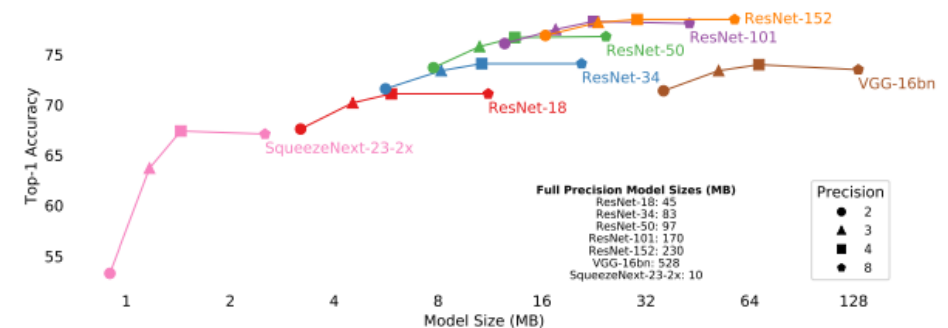
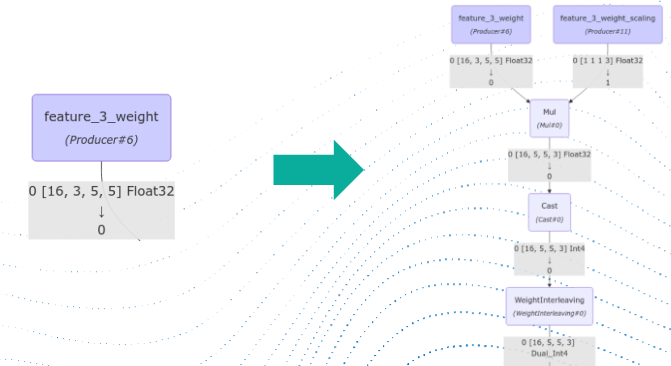




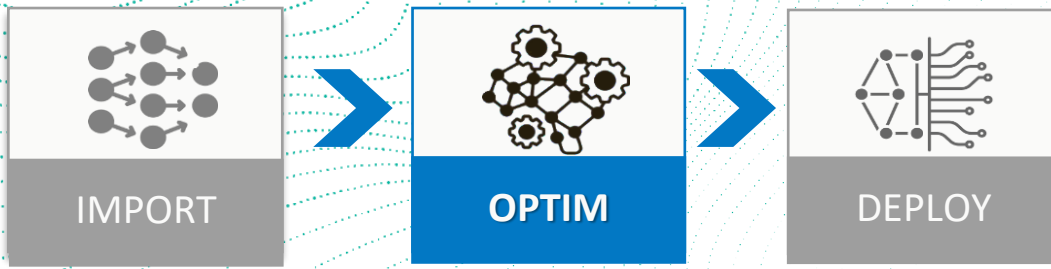
aidge

- **Automatic model reduction:** catalog of optimizations with deletion, reorganization and merging of operations
- **State-of-the-art quantification to desired accuracy (ResNet, VGG, etc.)**
 - After learning: without loss up to 8-bit integer
 - During learning: without loss up to 4-bit integer
- **Tensor decomposition compression method**
 - ResNet-50 x ImageNet: 15% compression without loss
 - ResNet-18 x CIFAR 100: 45% compression without loss

Traceability of optimization for certification purposes

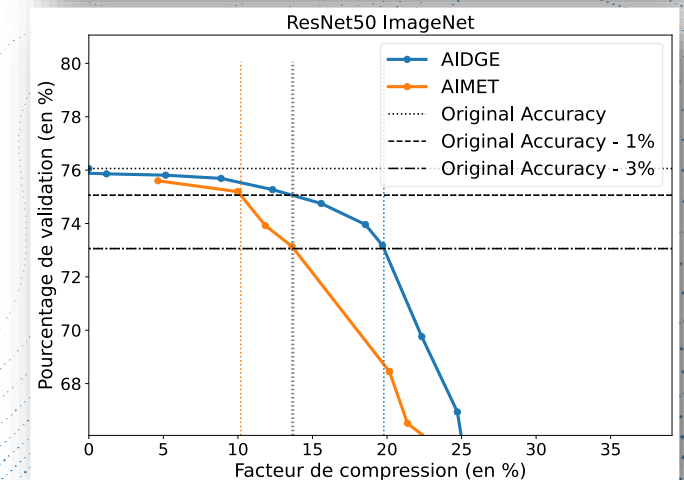
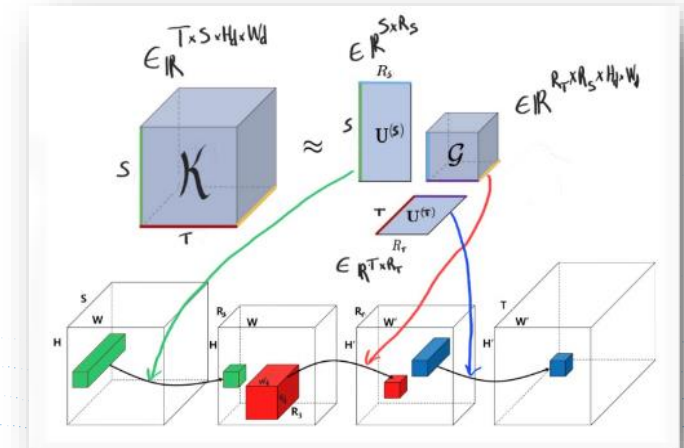


LEARNED STEP SIZE QUANTIZATION, Esser et al.

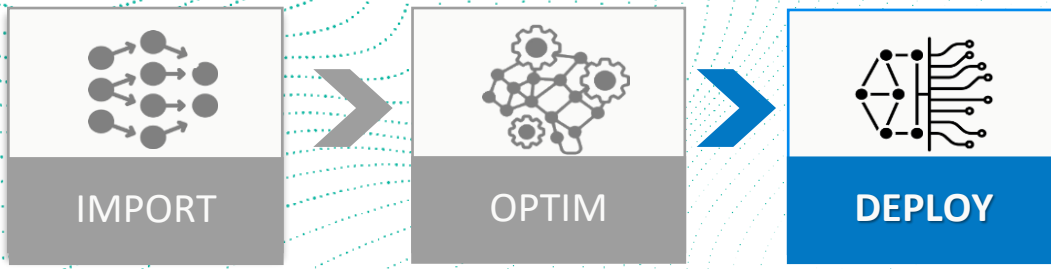


aidge

- **Automatic model reduction:** catalog of optimizations with deletion, reorganization and merging of operations
- **State-of-the-art quantification to desired accuracy (ResNet, VGG, etc.)**
 - After learning: without loss up to 8-bit integer
 - During learning: without loss up to 4-bit integer
- **Tensor decomposition compression method**
 - **ResNet-50 x ImageNet: 15% compression without loss**
 - **ResNet-18 x CIFAR 100: 45% compression without loss**



Higher performance than AIMET (Qualcomm)



aidge

- **ONNX export for interfacing with numerous SDKs**
- **Transparent, multi-paradigm code generation engine** (C/C++, HDL, etc.), enabling integration of compute kernels (native or third-party)
- **Multi-target reference export (C++) and specializations** (ARM, Texas Instrument SoC and ESP32 coming soon)
- **Orchestration control and memory optimization** through statistical allocation

STM32
Cube.AI

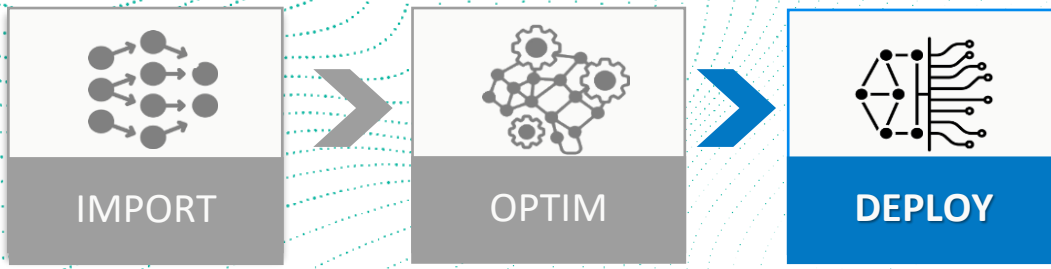
KaNN™
Kalray Neural Network

TensorRT

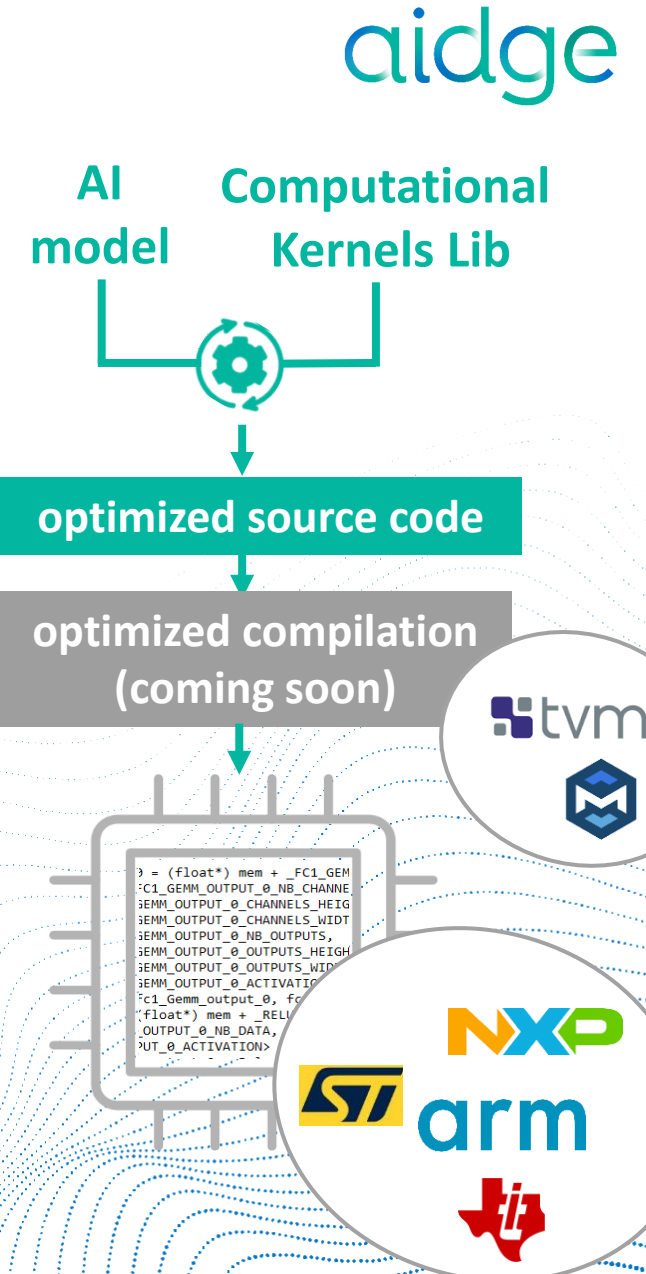
OpenVINO™

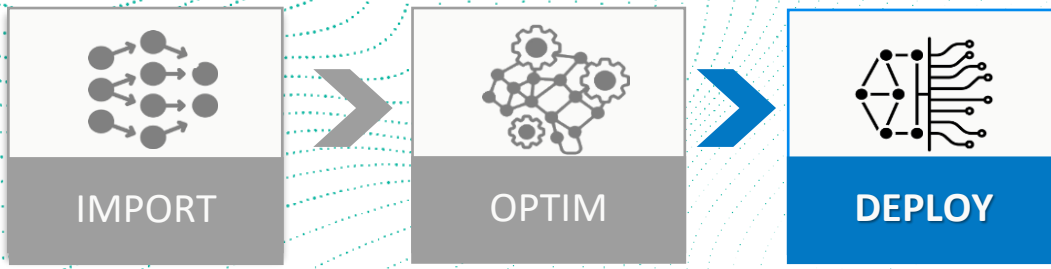
ONNX
RUNTIME

Qualcomm



- **ONNX export** for interfacing with numerous SDKs
- **Transparent, multi-paradigm code generation engine** (C/C++, HDL, etc.), enabling integration of compute kernels (native or third-party)
- **Multi-target reference export (C++)** and **specializations** (ARM, Texas Instrument SoC and ESP32 coming soon)
Coming soon :
 - Certification-aware export in C (ONERA)
 - Compilation workflow (INRIA)
- **Orchestration control and memory optimization** through statistical allocation

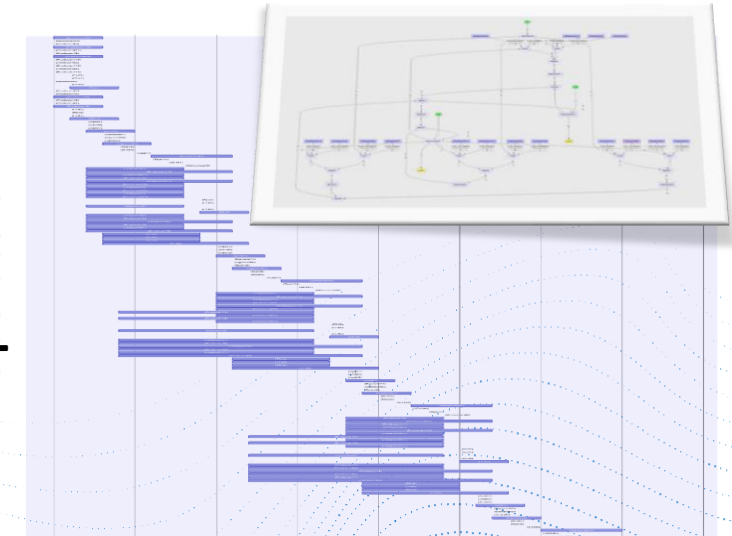




aidge

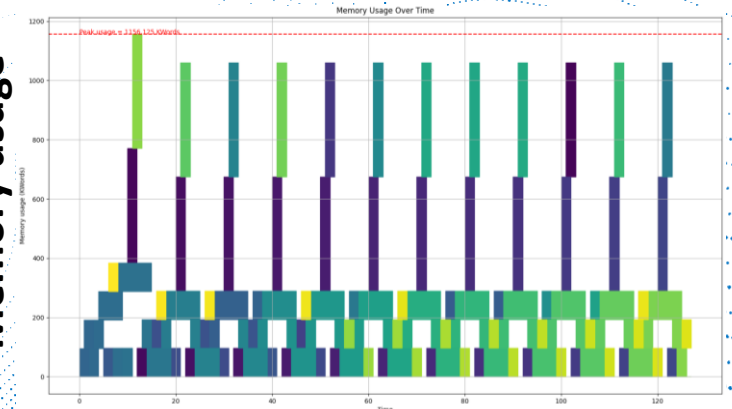
- **ONNX export** for interfacing with numerous SDKs
- **Transparent, multi-paradigm code generation engine** (C/C++, HDL, etc.), enabling integration of compute kernels (native or third-party)
- **Multi-target reference export (C++)** and specializations (ARM, Texas Instrument SoC and ESP32 coming soon)
- **Orchestration control and memory optimization** through statistical allocation

Operators



Time

Memory usage



Active developments and collaborations

+40

submitters

5

organizations

+300

Commits
/month

6th

version

20cent ikucher bhalimi obichler fabricea mnewson gkubler wboussella jeromeh raphaelmillet pineapple diegob marwaabd louislerbourg sylvainbataille alalloyer macario yberkat vbaudelet flebert mszczep farnez oantoni jsimatic nvrlosemyself idealbuquerque clementgf julienl nthm bobot alemesle thibaultallenet jgirardsatabin lucaslopez vlorrain axelfarr cguillon silvanosky hleborgne operrin mick94 na25 lsoulier cmoineau noamzerah alicebatte hrouis

Norms



afnor



ONNX

+30 industrials partners

THALES
Building a future we can all trust



AIRBUS



ALSTOM

MBDA
MISSILE SYSTEMS



ArcelorMittal

NX
NanoXplore



SAFRAN

KNDS

...

+10 academic partners



Fraunhofer



Inria

ONERA
THE FRENCH AEROSPACE LAB




université
PARIS-SACLAY

...


Some use cases

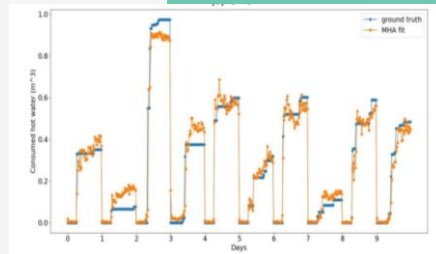
Defect detection and classification

- Low latency (20m/s) and high performance algorithm to detect small defect (~mm) with low contrast
- Deployment on Nvidia GPU
- In collaboration with 





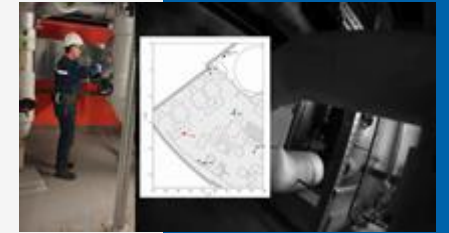
Heat Pump Monitoring

- Lightweight prediction algorithm based on incremental learning for adaptive heat pump control and monitoring
- Deployment on STM32
- Up to 40% energy saving
- In collaboration with 



Indoor localisation using multi-sensors

- Lightweight prediction algorithm based on IMU sensors combined with fast AI Visual tracking (x15 faster)
- Deployment on STM32
- In collaboration with  



Hardware design : NeuroCorgi AI-ASIC accelerator

- RTL generation of quantized model
- HD images processing in real time : latency is less than 10ms
- Uses 1,000 times less power than commercial circuits



aidge

Follow us



Join us and chat !

