

# Frameworks

# MindSpore (Huawei)

Prof. Carlos Moraes (Carlão)



**UNIFEI**



**Softex**



FUTURO DO TRABALHO, TRABALHO DO FUTURO

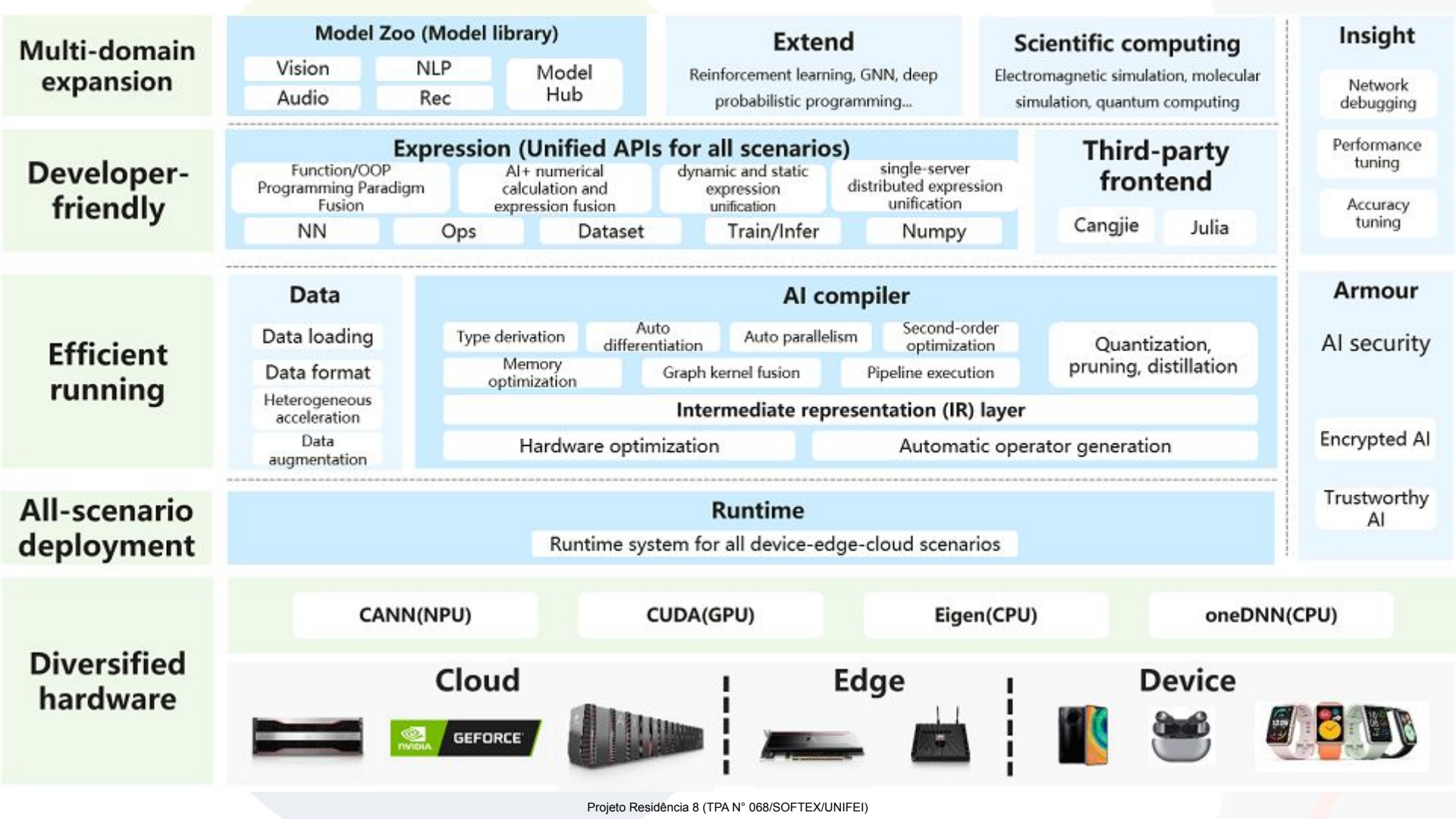


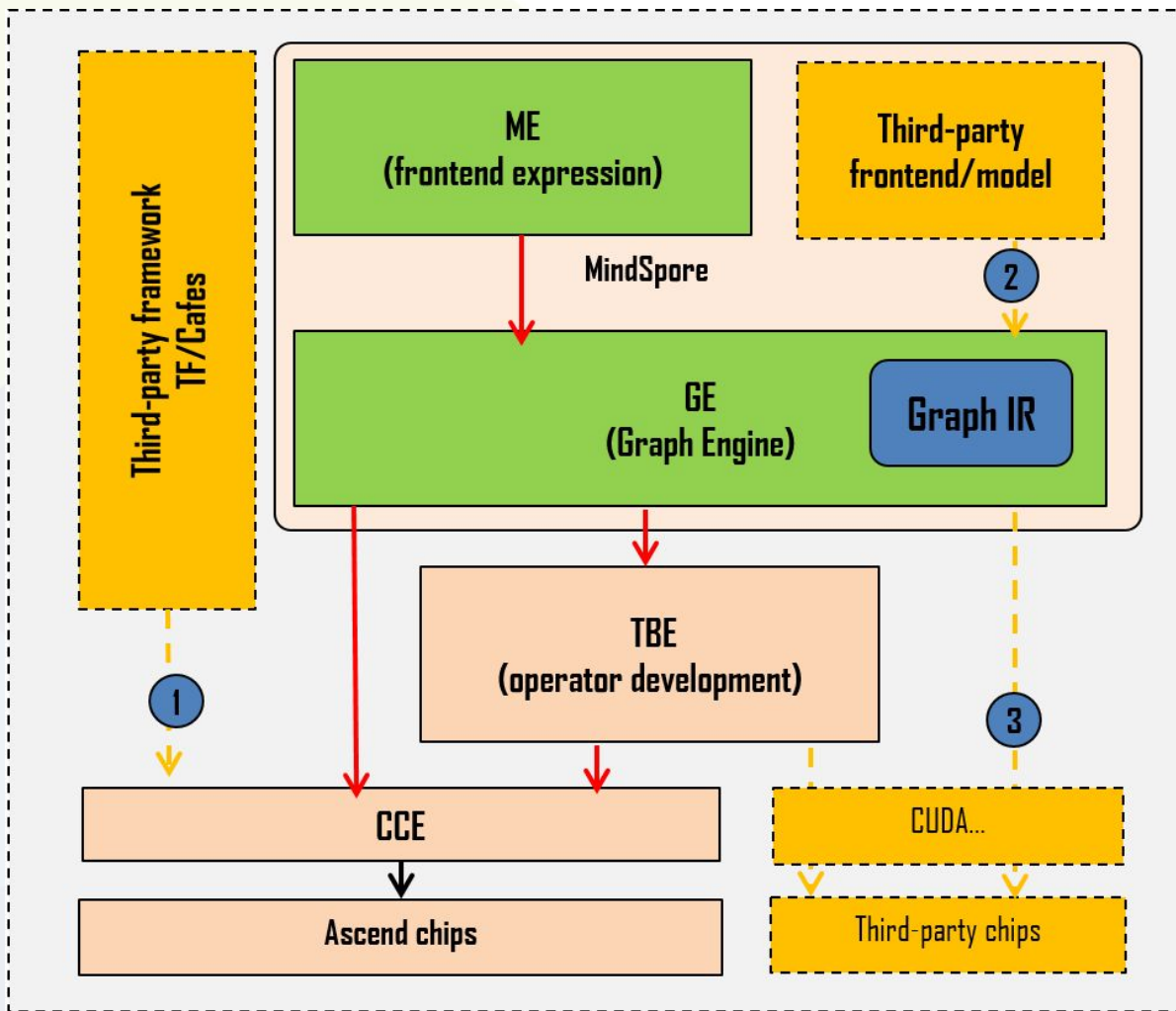
UNIÃO E RECONSTRUÇÃO



# [M]<sup>s</sup>

## MindSpore





# MindSpore

Unified APIs for all scenarios

Auto differ

Auto parallelism

Auto tuning

MindSpore intermediate representation (IR) for computational graph

On-device execution

Pipeline parallelism

Deep graph optimization

Device-edge-cloud co-distributed architecture (deployment, scheduling, communications, etc.)

Easy development:

AI Algorithm As Code

Efficient execution:

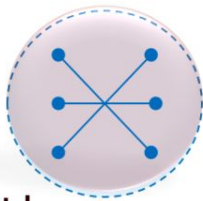
Optimized for Ascend

GPU support

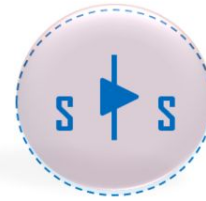
Flexible deployment: on-demand cooperation across all scenarios

**Processors: Ascend, GPU, and CPU**

# Diferenciação Automática



Technical path  
of automatic differential



## Graph: TensorFlow

- Non-Python programming based on graphs
- Complex representation of control flows and higher-order derivatives

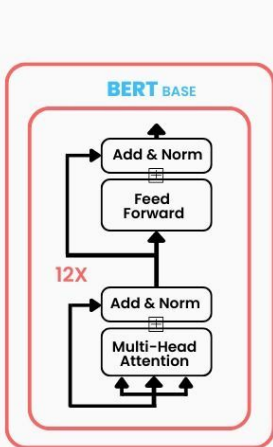
## Operator overloading: PyTorch

- Runtime overhead
- Backward process performance is difficult to optimize.

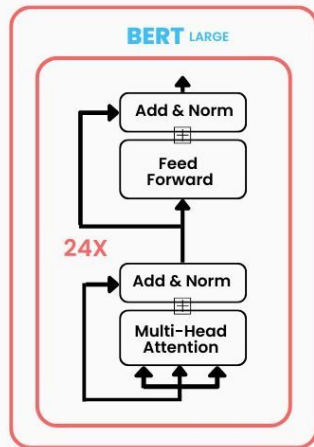
## Source code transfer: MindSpore

- Python APIs for higher efficiency
- IR-based compilation optimization for better performance

# Paralelismo Automático - Desafios



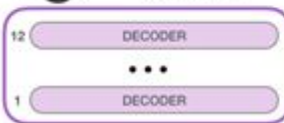
110M Parameters



340M Parameters



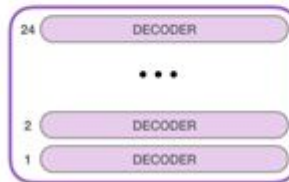
GPT-2  
SMALL



Model Dimensionality: 768



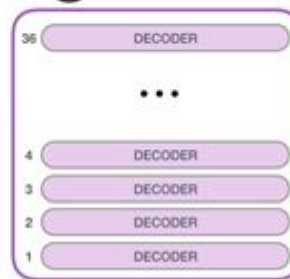
GPT-2  
MEDIUM



Model Dimensionality: 1024



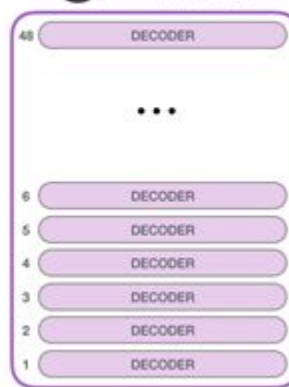
GPT-2  
LARGE



Model Dimensionality: 1280



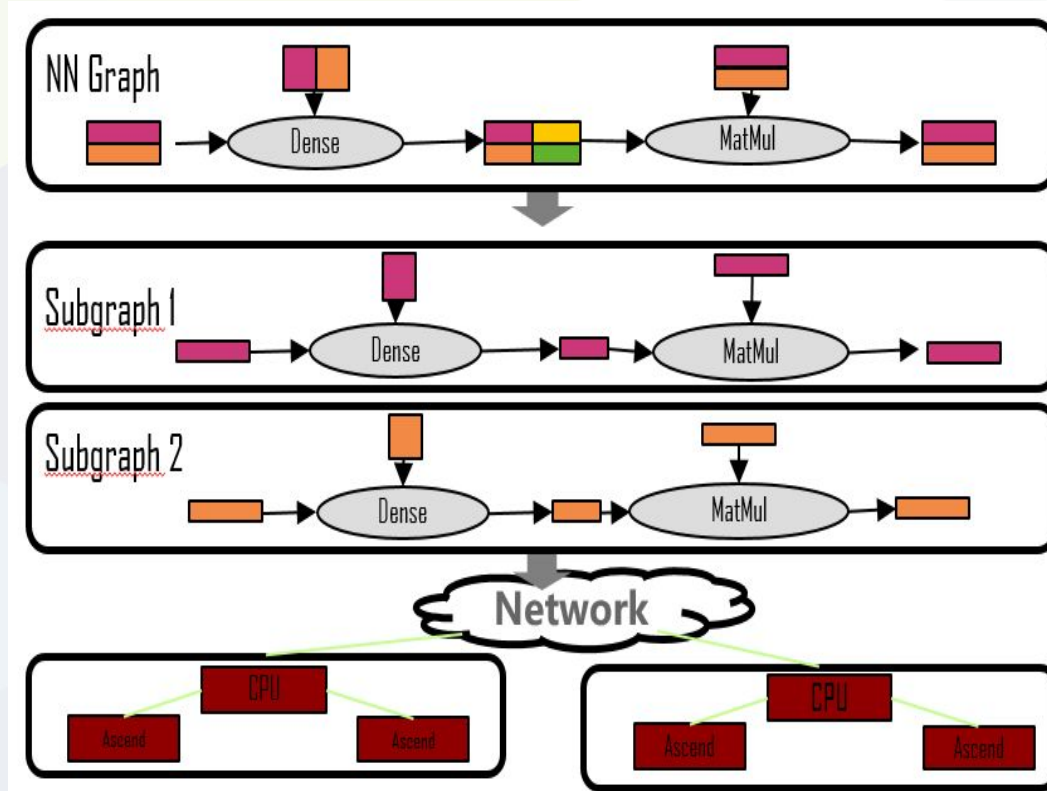
GPT-2  
EXTRA  
LARGE



Model Dimensionality: 1600

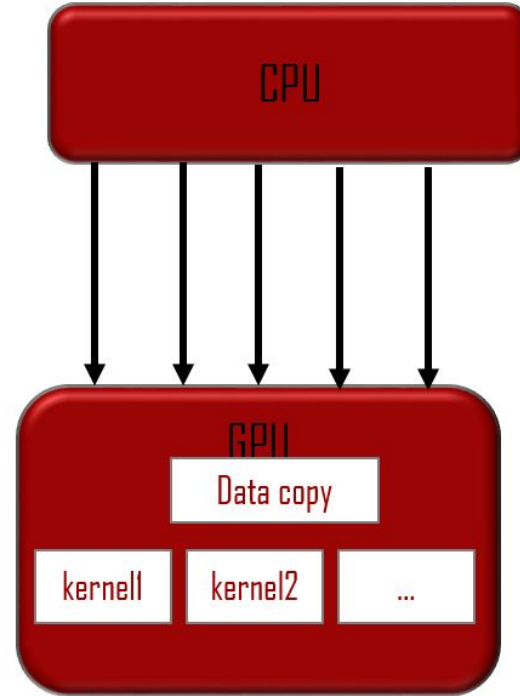
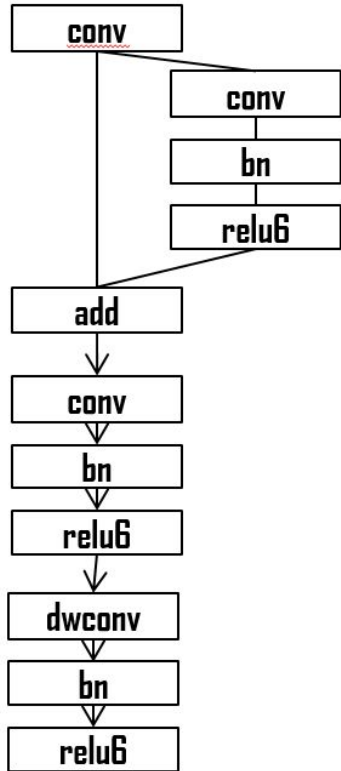


# Paralelismo Automático - Tecnologias



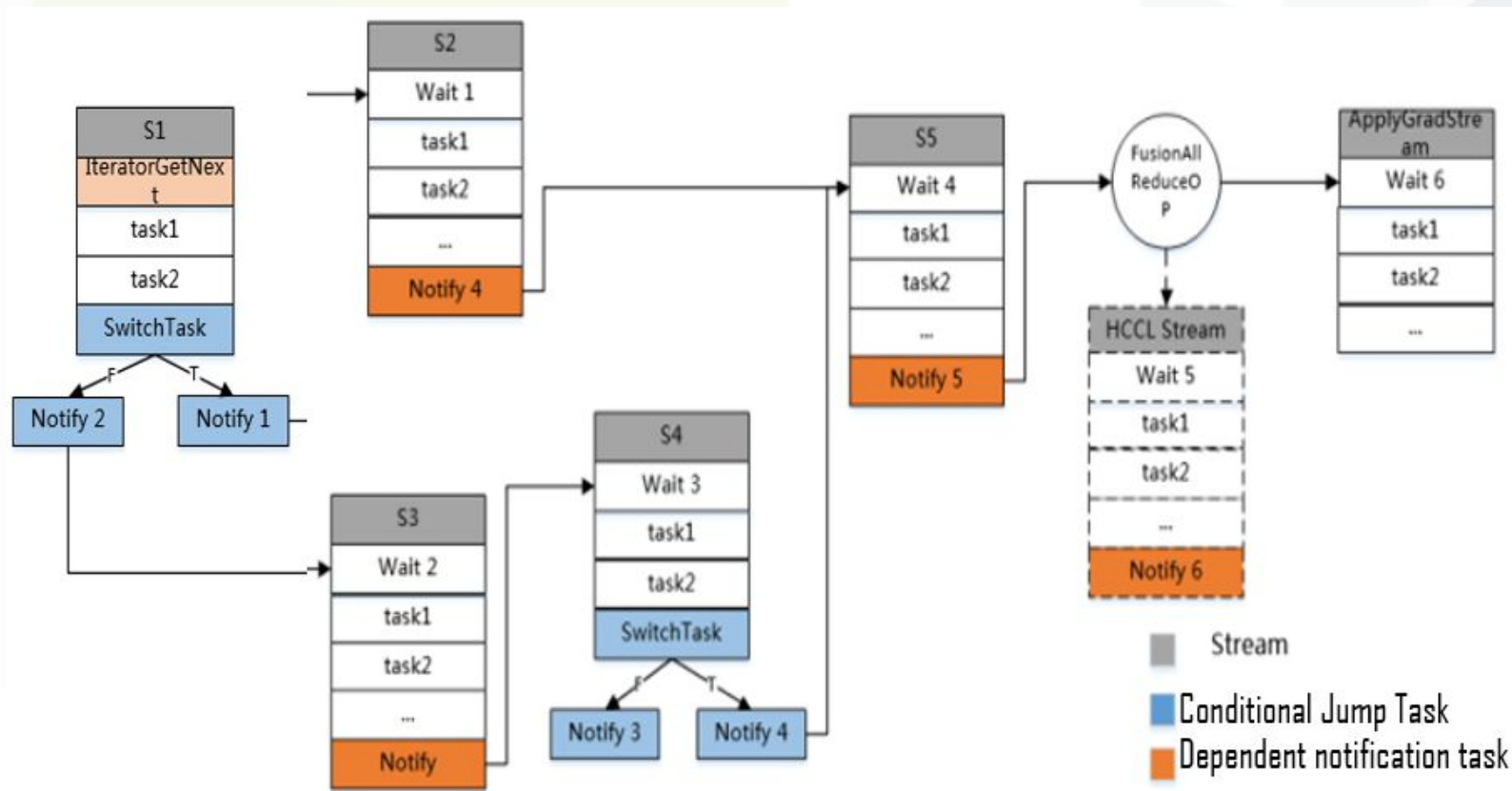


# Execução no Dispositivo

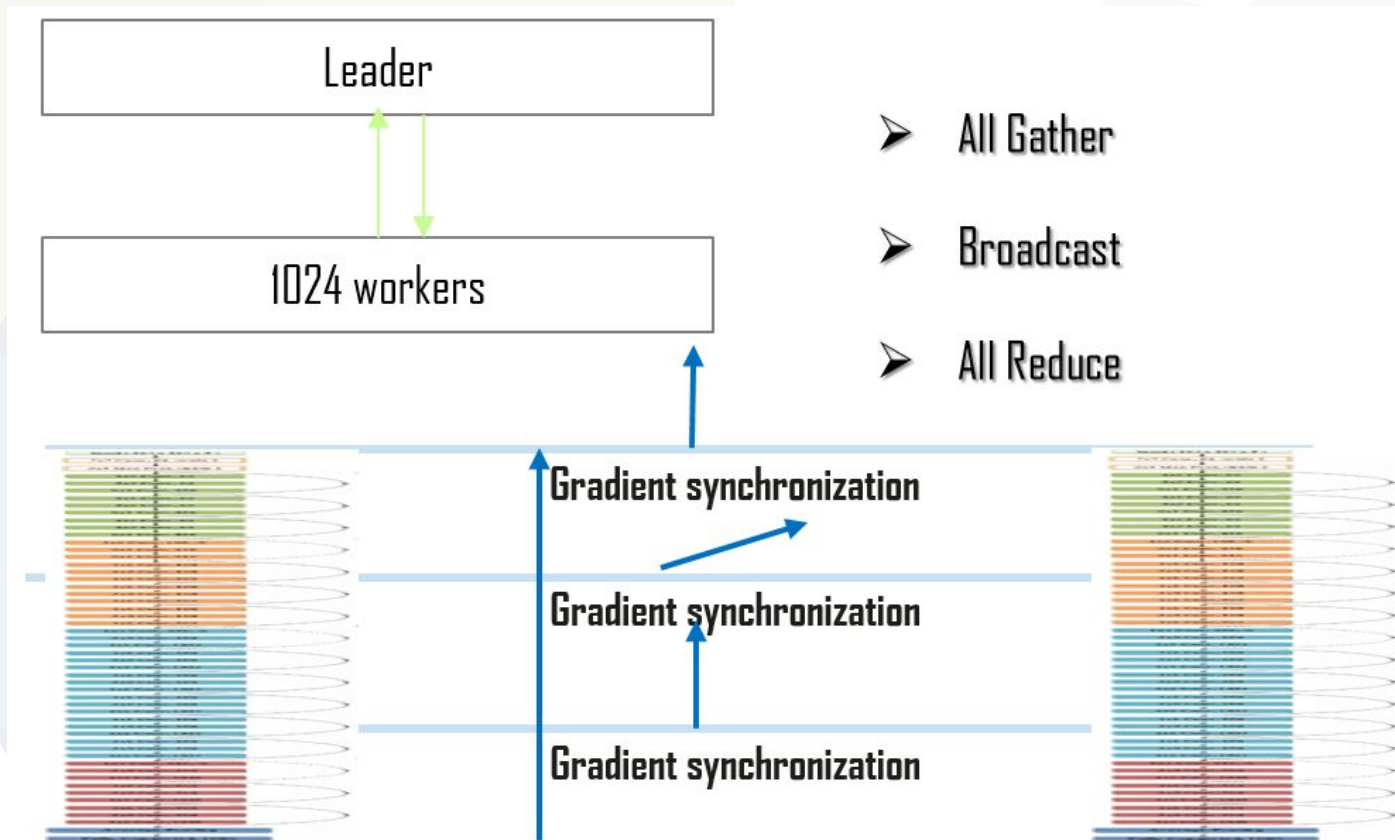


Large data interaction overhead and difficult data supply

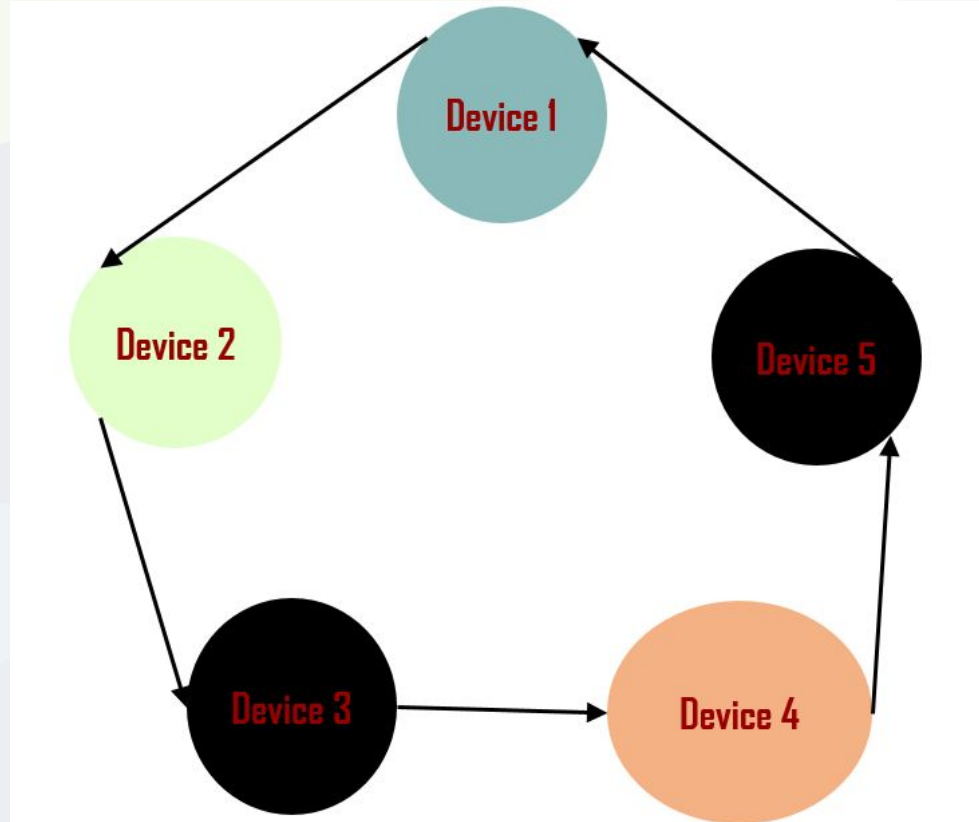
# Execução no Dispositivo - Tecnologias



# Desafios para Agregação de Gradiente Distribuído



# Tecnologias para Agregação de Gradiente Distribuído



# Arquitetura Distribuída Device-Edge-Cloud

On-demand collaboration in all scenarios and consistent development experience





# [M]<sup>s</sup>

## Características

# AI Computing Framework: Desafios

Desafios da Indústria

**Uma enorme lacuna entre pesquisa da indústria e aplicação de IA em todos os cenários**

- Barreiras de entrada altas
- Alto custo de execução
- Longa duração de implantação

Inovação tecnológica

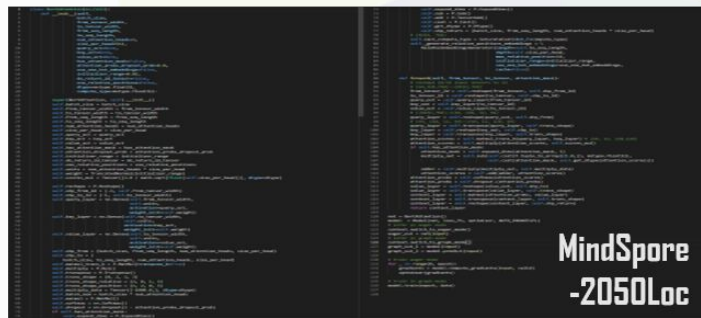
**MindSpore facilita IA inclusiva em todas as aplicações**

- Novo modo de programação
- Novo modo de execução
- Novo modo de colaboração

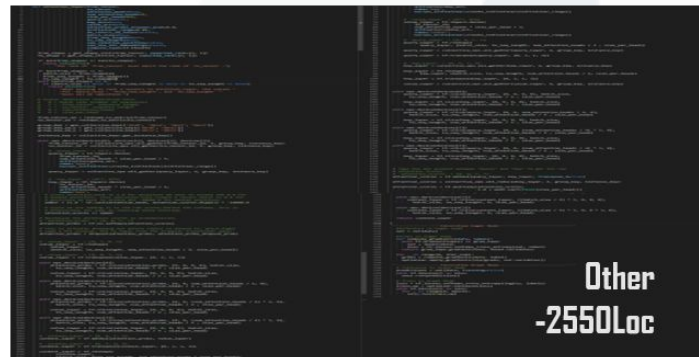


# Novo Paradigma de Programação

Algorithm scientist



Algorithm scientist  
+  
Experienced system developer



One-line  
Automatic parallelism

Efficient automatic differential

One-line debug-mode switch

NLP Model: Transformer

# TensorFlow vs MindSpore

```
1 import tensorflow as tf
2
3 model() {
4     with tf.device("/device:0")
5         token_type_table = tf.get_variable(
6             name=token_type_embedding_name,
7             shape=[token_type_vocab_size, width],
8             initializer=create_initializer(initializer_range))
9         flat_token_type_ids = tf.reshape(token_type_ids, [-1])
10        one_hot_ids = tf.one_hot(flat_token_type_ids, depth=token_type_vocab_size)
11        token_type_embeddings = tf.matmul(one_hot_ids, token_type_table)
12
13    with tf.device("/device:1")
14        query_layer = tf.layers.dense(
15            from_tensor_2d,
16            num_attention_heads * size_per_head,
17            activation=query_act,
18            name="query",
19            kernel_initializer=create_initializer(initializer_range))
20
21    with tf.device("/device:2")
22        key_layer = tf.layers.dense(
23            to_tensor_2d,
24            num_attention_heads * size_per_head,
25            activation=key_act,
26            name="key",
27            kernel_initializer=create_initializer(initializer_range))
```

```
class DenseMatMulNet(nn.Cell):
    def __init__(self):
        super(DenseMatMulNet, self).__init__()
        self.matmul1 = ops.MatMul.set_strategy([4, 1], [1, 1])
        self.matmul2 = ops.MatMul.set_strategy([1, 1], [1, 4])

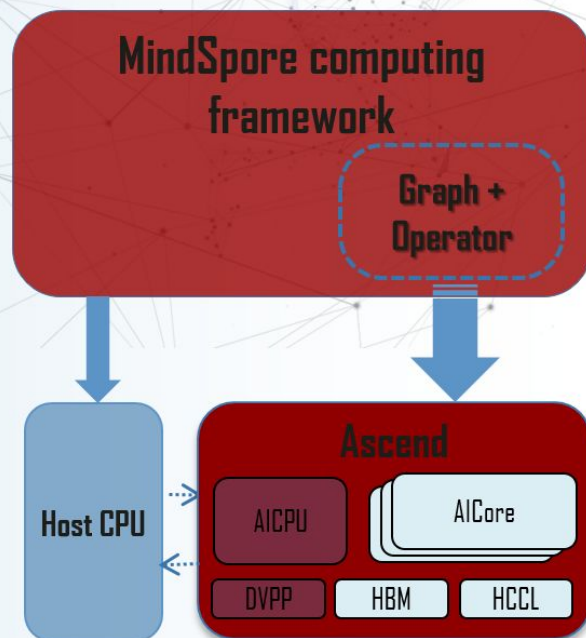
    def construct(self, x, w, v):
        y = self.matmul1(x, w)
        z = self.matmul2(y, v)
        return z
```

# Desafios de Execução

## On-device execution

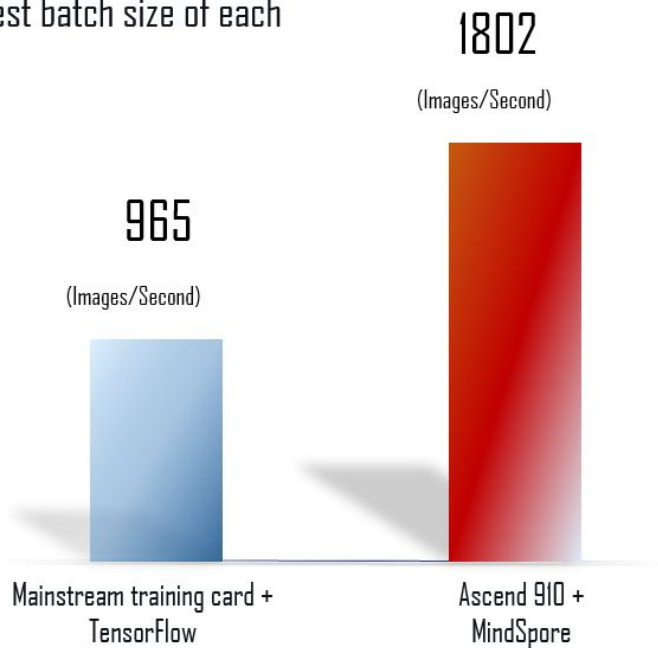
Offloads graphs to devices, maximizing the computing power of Ascend

- ◆ **Framework optimization**
  - Pipeline parallelism
  - Cross-layer memory overcommitment
- ◆ **Soft/hardware co-optimization**
  - On-device execution
  - Deep graph optimization



# O desempenho do ResNet-50 é dobrado

- ResNet 50 V1.5
- ImageNet 2012
- With the best batch size of each card



# Reconhecimento de vários objetos em tempo real



**Detecting objects in 60 ms**

# Novo modo de colaboração



**V.S.**



# Diferentes precisão e velocidade de hardware





# Desenvolvimento unificado e Implantação flexível



# Alta performance

## AI computing challenges

### Complex computing

Scalar, vector, and tensor computing

Mixed precision computing

Parallelism between gradient aggregation  
and mini-batch computing

### Diverse computing units / processors

CPUs, GPUs, and Ascend processors

Scalar, vector, and tensor

## MindSpore

### Framework optimization

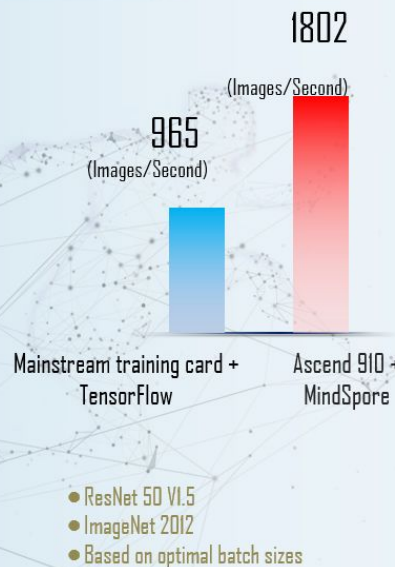
Pipeline parallelism

Cross-layer memory overcommitment

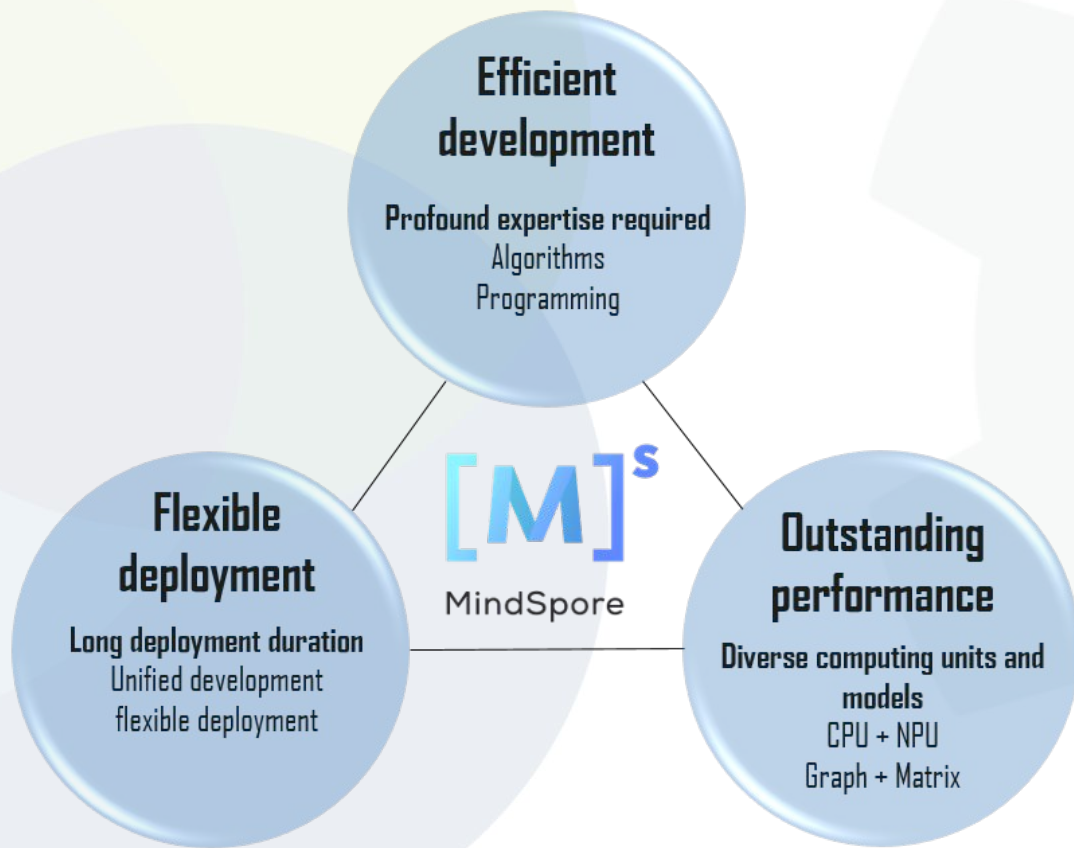
### Software + hardware co-optimization

On-device execution

Deep graph optimization



# Visão e valor da MindSpore

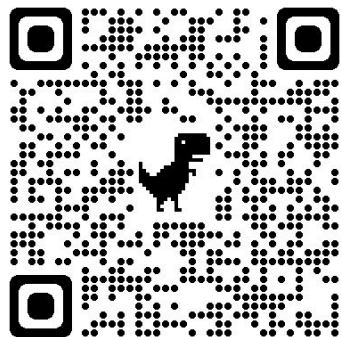




# [M]<sup>s</sup>

## Desenvolvimento e Aplicação

# Instalando o MindSpore



← → ↻ mindspore.cn/install/en

[M] MindSpore Install Tutorials Docs Ecosystem Community Resources News Security

This site uses cookies. By continuing to browse the site you are agreeing to our use of cookies. [Read our privacy policy](#)

Select an environment, obtain commands, and perform installation according to the instructions. Alternatively, create and deploy a model on ModelArts.

### Obtaining Installation Commands

[View All Versions and API Updates](#)

Version	2.0.0-alpha	1.10.0 ✓	Nightly		
Hardware Platform	Ascend 910	Ascend 310	GPU CUDA 10.1 ✓	GPU CUDA 11.1	CPU
Operating System	Linux-aarch64	Linux-x86_64 ✓	Windows-x64	MacOS-aarch64	MacOS-x86_64
Programming Language	Python 3.7 ✓	Python 3.8	Python 3.9		
Installation Mode	Pip ✓	Conda	Source	Docker	Binary
Commands	<pre>pip install https://ms-release.obs.cn-north-4.myhuaweicloud.com/1.10.0/MindSpore/gpu/x86_64/cuda-10.1/mindspore_gpu-1.10.0-cp37-cp37m-linux_x86_64.whl --trusted-host ms-release.obs.cn-north-4.myhuaweicloud.com -i https://py.pi.tuna.tsinghua.edu.cn/simple</pre> <p># Refer to the following installation guide, ensure that the installation dependency and environment variables are correctly configured.</p>				

Módulos

**model\_zoo**

Define modelos de rede comuns

**communication**

Módulo de carregamento de dados, que define o dataloader e dataset e processa dados como imagens e textos.

**dataset**

Módulo de processamento de dataset, que lê e pré-processar dados.

**common**

Define tensor, parâmetro, dtype e inicializador.

**context**

Define a classe de contexto e define os parâmetros de execução do modelo, como grafos e PyNative modos de comutação.

# Módulos

**akg**

Diferencial automático e biblioteca de operador personalizado.

**nn**

Define células MindSpore (unidades de rede neural), funções de perda e otimizadores.

**ops**

Define operadores básicos e registra operadores reversos.

**train**

Modelo de treinamento e módulos de função de resumo.

**utils**

Utilitários, que verificam os parâmetros. Este parâmetro é usado na estrutura.

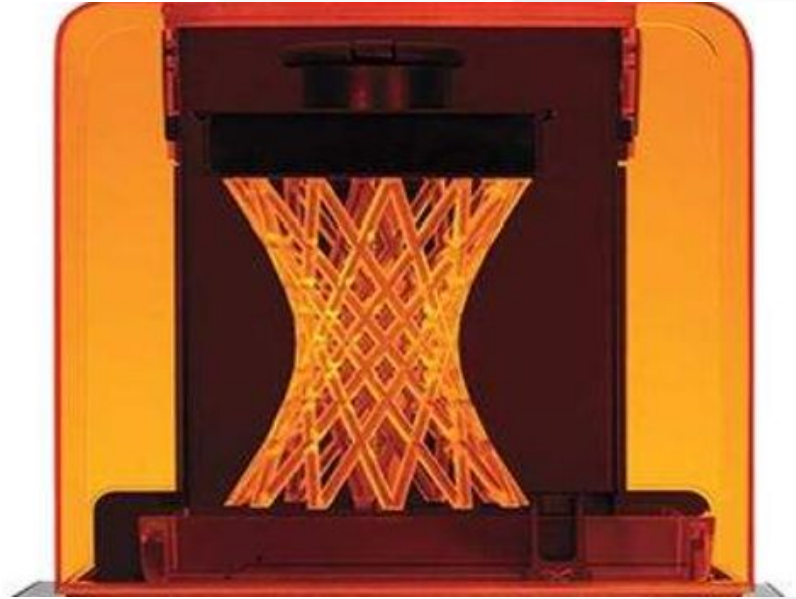
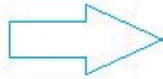


# Operações com Tensores

- `asnumpy()`
- `size()`
- `dim ()`
- `dtype()`
- `set_dtype()`
- `tensor_add(Tensor)`
- `tensor_mul(Tensor)`
- `shape()`
- `__Str__ #` (conversão em strings)

# Célula MindSpore

```
class Net(nn.Cell):  
    def __init__():  
        .....  
    def construct(self, x):  
        .....
```



# Pergunta

Which of the following is a math operator in MindSpore?

- Mul
- Select
- Callback
- ControlDepend

# Pergunta

Automatic graph segmentation is one of the key technologies of MindSpore.

- Operator-level automatic differentiation
- Automatic parallelization
- Automatic operator generation
- Semi-automatic data labeling

# Pergunta

Facing industry research and full scenarios AI The huge gap between applications, MindSpore Bridging the application gap to help inclusiveness AI of technological innovation does not include which of the following?

- New programming language
- New ways of collaboration
- New programming paradigm
- New execution mode

# Pergunta

In MindSpore, if you want to define a function for collecting information after each step is complete, you need to inherit the `[?]` class.

# Pergunta

The [?] subsystem of MindSpore provides Mind IR-oriented, graph-level, real-time compilation.



# Apoio

Este projeto é apoiado pelo Ministério da Ciência, Tecnologia e Inovações, com recursos da Lei nº 8.248, de 23 de outubro de 1991, no âmbito do [PPI-Softex| PNM-Design], coordenado pela Softex.



**UNIFEI**



**Softex**



PARCELO 01 - TRILHA DO FUTURO





# Softex

**MCTI  
FUTURO**