# Guillaume Moigneu

VP, Advocacy | **Upsun & Platform.sh**

**nls.io on bluesky**

# Qui est familier avec les Transformers?

# Qu'est-ce qu'un Transformer?

Une architecture de réseau neuronal qui a révolutionné le NLP et au-delà.

Introduite dans l'article de 2017 "Attention is All You Need"
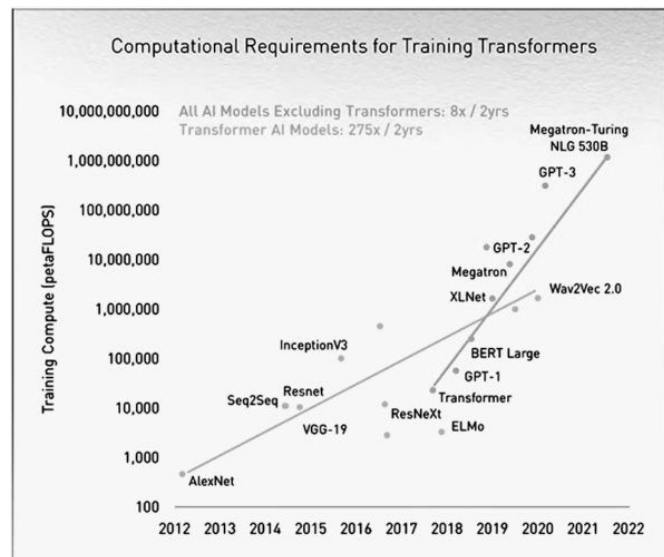Alimente des modèles comme GPT, BERT et T5

**Une façon de gérer une tâche spécifique au-dessus d'un modèle. Plus efficacement.**
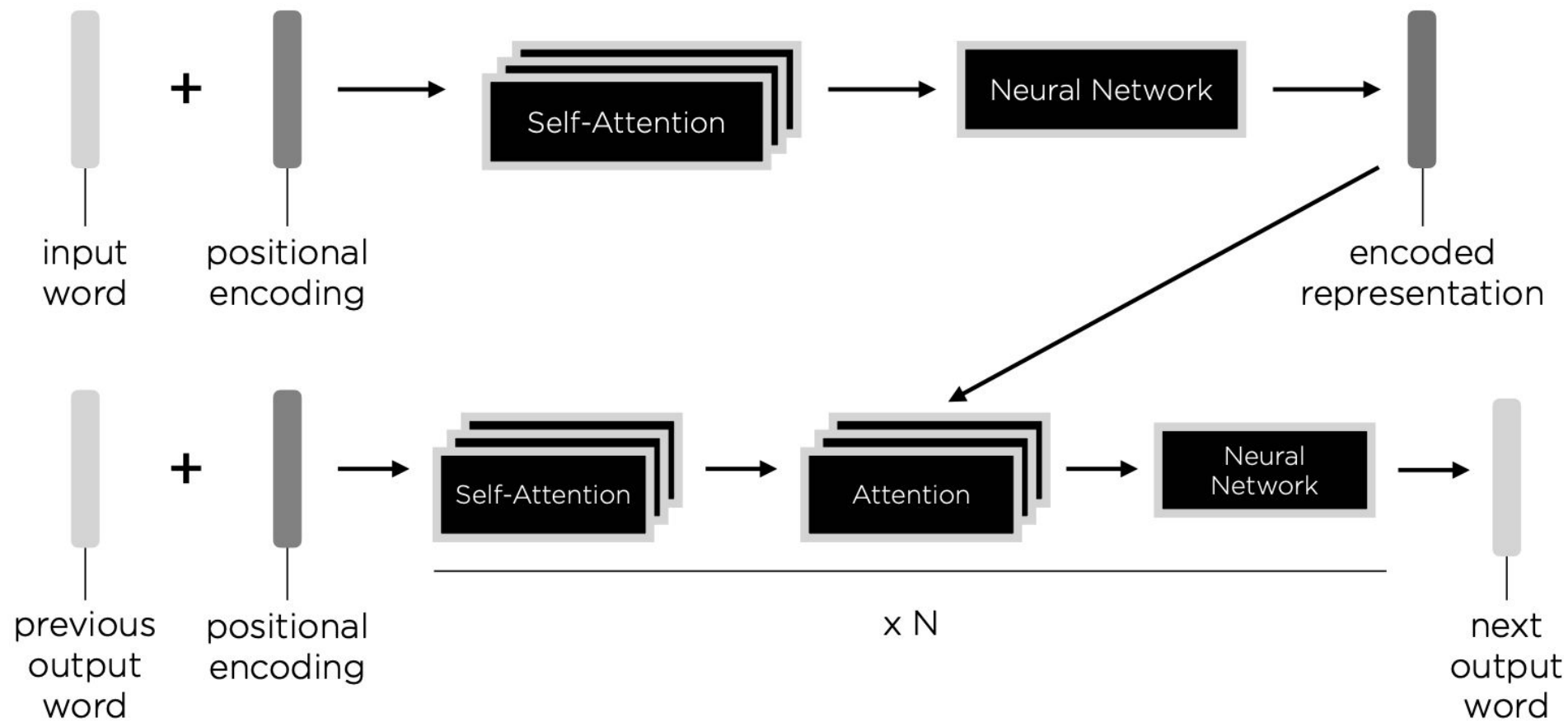
# Innovation clé : Traite toutes les entrées simultanément au lieu de séquentiellement

Structure encodeur-décodeur
Couches d'auto-attention
Réseaux de neurones feed-forward
Encodages de position

**Avantages:**

- Traitement parallèle et réduction du coût d'entraînement et d'inférence
- Compréhension étendue du contexte
- Entraînement et inférence plus rapides
- Meilleur captage des relations
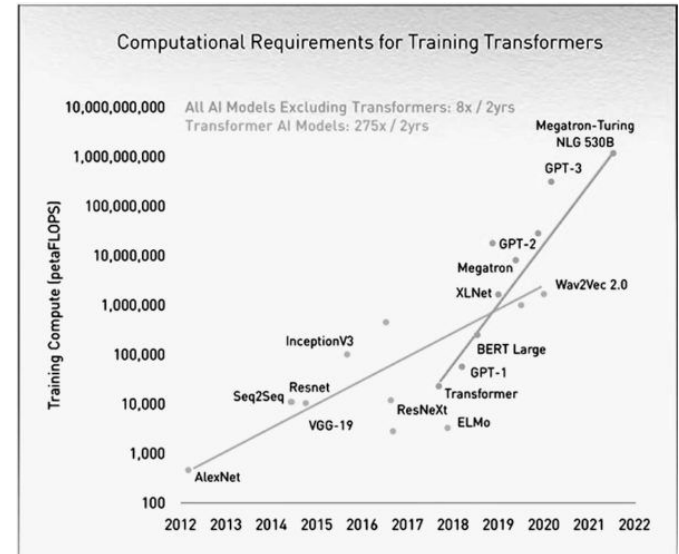- Utilisation réduite des ressources



Computational Requirements for Training Transformers

All AI Models Excluding Transformers: 8x / 2yrs
Transformer AI Models: 275x / 2yrs

input word + positional encoding → Self-Attention → Neural Network → encoded representation

previous output word + positional encoding → Self-Attention → Attention → Neural Network → next output word

x N

# Démo et liens



**Demo:**

https://huggingface.co/spaces/webml-community/attention-visualization

**Learn more:**

https://cdn.cs50.net/ai/2023/x/lectures/6/lecture6.pdf

Inférence Machine Learning en PHP par l'exemple

# Implémentation de 3 cas d'usages en PHP

# Pourquoi faire de l'inférence locale en PHP?

- Aucune souscription tierce
- Pas de "fuite" de données à l'extérieur
- Pas de changement d'architecture/setup

- Millions de modèles disponibles + vos propres modèles
- Fine-tuning possible
- Meilleur contrôle (test A/B, performance, uptime)



Computational Requirements for Training Transformers

All AI Models Excluding Transformers: 8x / 2yrs
Transformer AI Models: 275x / 2yrs

```
composer require codewithkyrian/transformers # Package

./vendor/bin/transformers install # Install platform specific transformers

# Download models
./vendor/bin/transformers download Xenova/distilbert-base-uncased-finetuned-sst-2-english # Text
classification
./vendor/bin/transformers download Xenova/vit-base-patch16-224 # Image classification
```

# Classification de texte

```php
use function Codewithkyrian\Transformers\Pipelines\pipeline;

protected function execute(InputInterface $input, OutputInterface $output): int
{
    $asin = $input->getArgument('asin');
    $output->writeln("Calculating score for ASIN: $asin");

    $product = $this->entityManager->getRepository(Product::class)->findOneBy(['asin' => $asin]);
    $reviews = $product->getReviews();

    $pipe = pipeline('sentiment-analysis');

    foreach ($reviews as $review) {
        $out = $pipe($review->getText());
        if ($out['label'] == 'POSITIVE') {
            $positive++;
        } else {
            $negative++;
        }
    }

    $output->writeln("Positive: $positive, Negative: $negative, Score: " . ($positive - $negative) ."
(".round($positive / $reviews->count() * 100, 2)."% positive)");
    return Command::SUCCESS;
}
```

MM
FR

```
php bin/console app:score B07VGRJDFY
```

```
docker-compose          ~/psh/php-transformers                    Share

~/psh/php-transformers git:(main) ⌄ ±1    Pair ⌘ |    Dispatch Beta ⇧ ⌘ |
```

Cas d'usage #2

—

# Classification d'images

```php
public function handle(Request $request): Response
{
  [...]
  // Process the image to generate labels
  Transformers::setup()->setImageDriver(ImageDriver::GD);
  $classifier = pipeline('image-classification');
  $result = $classifier($this->getParameter('kernel.project_dir').'/public/uploads/'.$newFilename, 3);
  // Loop results to see if it's a hot dog
  [...]
}
```

https://mm25fr.moigneu.net/hotdog

# Génération de texte

```php
public function handle(Request $request): Response
{
    $question = $request->request->get('question');

    $generator = pipeline('text2text-generation', 'Xenova/flan-t5-small');
    $result = $generator($question,
        maxNewTokens: 256,
        repetitionPenalty: 1.6,
        temperature: 0.7
    );
    $answer = $result[0]['generated_text'];

    return $this->render('question/handle.html.twig', [
        'question' => $question,
        'answer' => $answer
    ]);
}
```

MM
FR

[https://mm25fr.moigneu.net/question](https://mm25fr.moigneu.net/question)

Inférence Machine Learning en PHP par l'exemple

Il est temps d'expliquer la magie

# PHP FFI - Foreign Function Interface

―――――

FFI extension allows PHP code to directly **call functions and manipulate data from C libraries** without writing additional C code or PHP extensions.

MM
FR

# Open Neural Network Exchange

___

ONNX (Open Neural Network Exchange) is an **open standard format** that allows AI models to be shared between different machine learning frameworks like PyTorch, TensorFlow, and many others.

# Transformers & pipelines

Fournis par le package TransformersPHP.

Managé par **Kyrian Obikwelu**
Et 7 contributeurs principaux

[Pipelines supportées](#)

# Models

Centaines de modèles dispos sur 🤗.

Vous pouvez convertir n'importe quel modèle en format ONNX via

- [Python directement](#)
- ou [un Notebook](#)



#MM25FR

# SmolLM

Modèle LLM prometteur et efficace

Fourni par l'équipe HF
Seulement 538 Mo
Excellente alternative aux modèles complets
Peut même fonctionner sur le navigateur

Inférence Machine Learning en PHP par l'exemple

Plus grande limitation actuellement :
Pas de support GPU. Pas encore.

## Ma recommandation

Bien que cela soit **excellent pour les petites tâches** (classification, étiquetage d'images, etc.), tout travail sur les LLM devrait être **délégué à une infrastructure basée sur du GPU**.

Utilisez un service SaaS : OpenAI, API Claude, etc.

Ou utilisez les points d'accès d'inférence HF pour déployer le modèle de votre choix
Et interrogez ce point d'accès API depuis votre application.

Ou "Construisez le vôtre", à vos risques et périls !

Meet **Magento** FRANCE

# MERCI !

Guillaume Moigneu
nls.io on bluesky

*plz follow!*