upsun

# Machine Learning Inference in PHP by example

Leverage ONNX and Transformers on Symfony
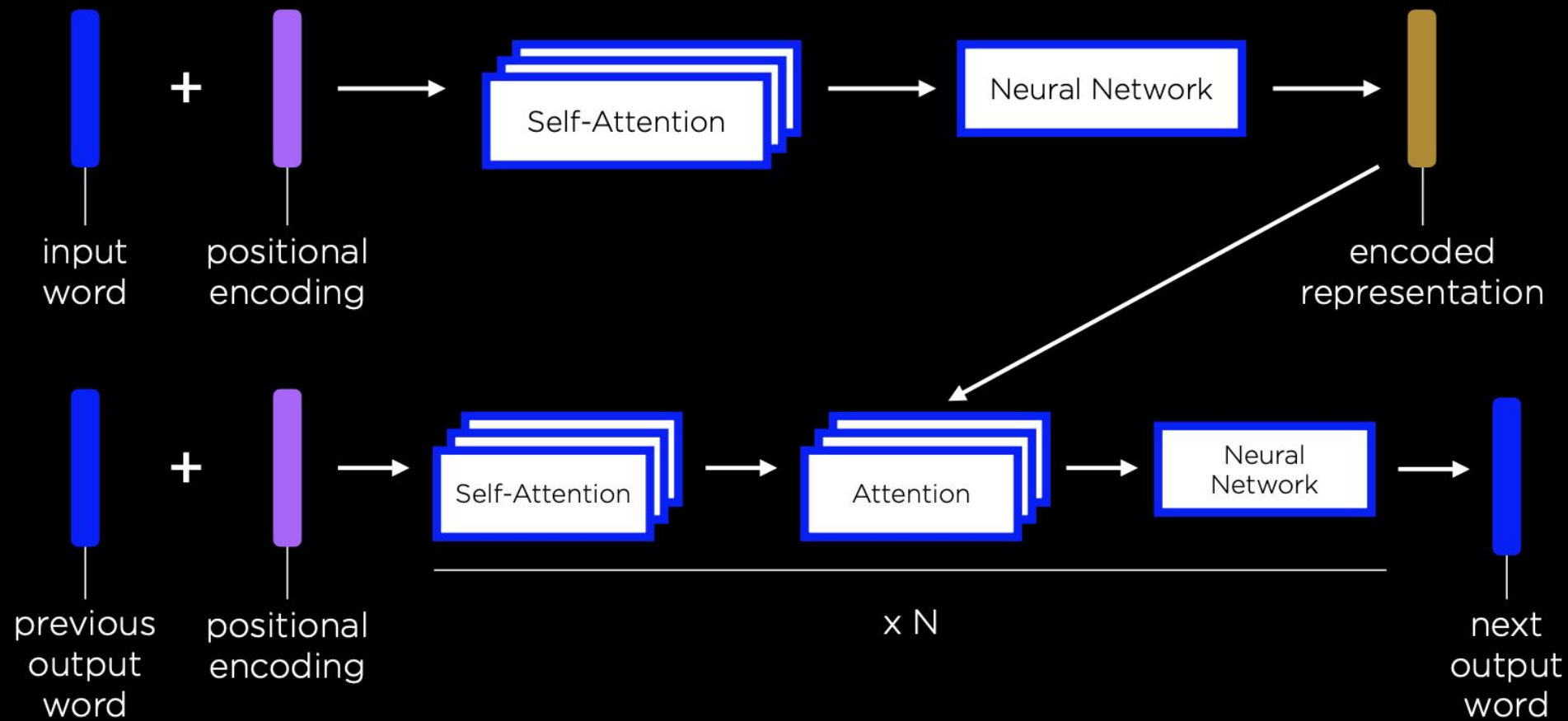
Migrate your app    Create from scratch    Migrate from GitHub

upsun environment:branch staging

hello, world!

Hello world

🌱 Canada

~28 gCO2eq/kWh

+ Team

Frontend developers

Production: Contributors
Staging: Contributors

Instant clone env
(+ all your data!)

dj  django-app

NEXT.  nextjs-app

database

production

branch

# Who is familiar with Transformers?

# What is a **Transformer?**

- A neural network architecture that revolutionized NLP and beyond

- Introduced in 2017 paper "Attention is All You Need"

- Powers models like GPT, BERT, and T5

## Key innovation: processes all input at once instead of sequentially

- Encoder-decoder structure
- Self-attention layers
- Feed-forward neural networks
- Position embeddings

# Before 2017 After

**Before**    **2017**    **After**

Sequential processing (RNNs, LSTMs)
Limited context window
Slow training and inference

Parallel processing & lowering cost of training & inference
Extended context understanding
Faster training and inference
Better at capturing relationships
Reduced resources usage

# Demo:

**Text: Symfony is a great framework for developing [MASK] applications.**

# Learn more:

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


SCAN ME

# Let's implement 3 use-cases
## in PHP

```
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
```

**Text Classification**

```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;
}
```

# Demo

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

Image Classification

```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
    [...]
}
```
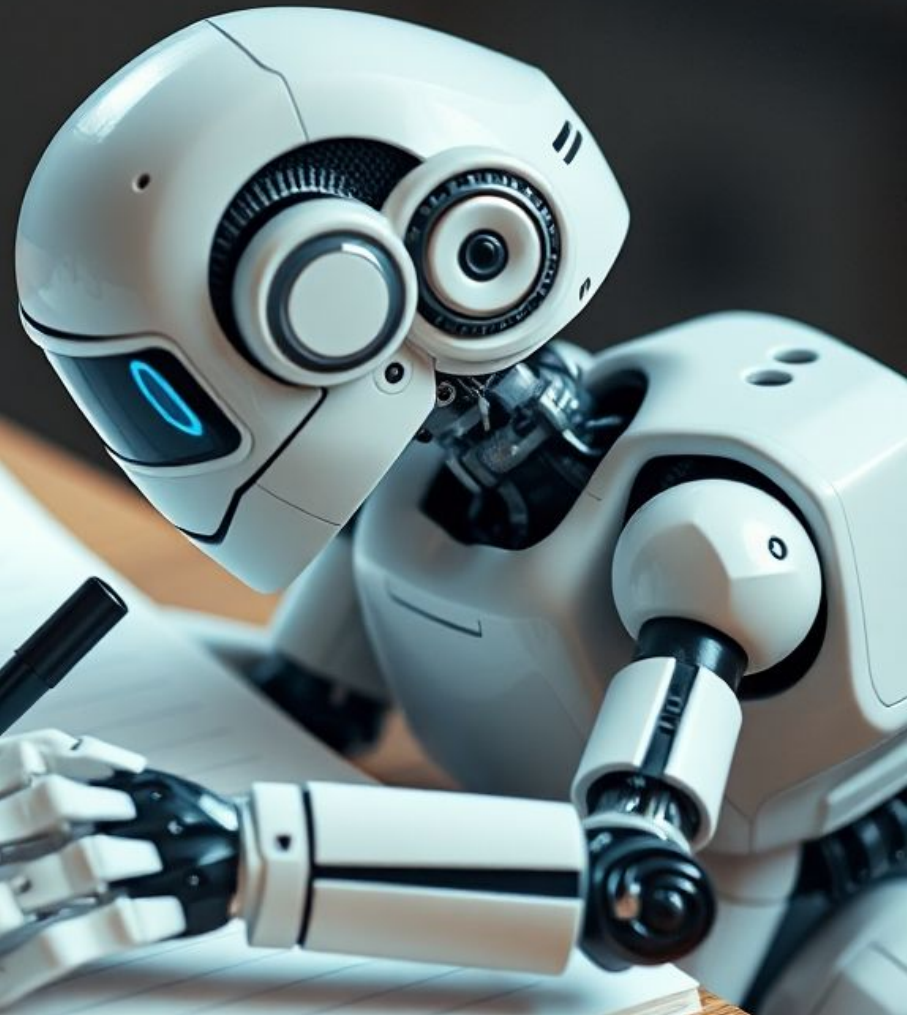
# Demo

[https://vienna.moigneu.net/hotdog](https://vienna.moigneu.net/hotdog)

LLM generation

```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
    ]);
}
```

# Demo

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

Time to explain the magic.

# 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.

# Open Neural Network Exchange

ONNX for short

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

and pipelines

Provided by the TransformersPHP package.

Courtesy of Kyrian Obikwelu

And 7 contributors

[Supported pipelines](#)

# Models

Infinite possibilities

- Hundreds of models readily available on 🤗.
- Convert any model to ONNX with:
  + [Python](#)
  + [Notebook](#)

# SmolLM

Promising efficient LLM model

- Provided by the HF team
- 538Mb only
- Great alternative to fully fledged models
- Can even run on the browser



🤗 Spaces · 🤗 HuggingFaceTB / **SmolLM-360M-Instruct-WebGPU** 🗂️    ♡ like    117    ● Running    🅰 App

👤    Give me some tips to improve my public speaking skills

🤖    Here are some tips to help you improve your public speaking skills:

1. **Practice, practice, practice**: The more you practice, the more comfortable you'll become with your ma[...] the more confident you'll feel. Practice in front of a mirror, record yourself, or speak in front of a small gr[...] friends or family members.

2. **Know your audience**: Understand who your audience is, what they're interested in, and what they're loo[...] Tailor your message, language, and tone to engage and persuade them.

3. **Prepare thoroughly**: Research your topic, organize your thoughts, and create an outline. Make sure you[...] clear structure and a solid understanding of your key points.

4. **Focus on your message**: Instead of worrying about how you're coming across, focus on the value you're[...] providing to your audience. Remember, you're speaking to share your ideas, not to impress anyone with[...] speaking skills.

5. **Make eye contact**: Engage with your audience by making eye contact with different people in the room.[...] helps build trust and keeps your audience interested.

6. **Use confident body language**: Stand up straight, make use of gestures, and avoid fidgeting or crossing[...] arms. This will help you appear more confident and others will perceive you as more confident too.

7. **Speak slowly and clearly**: Avoid rushing through your speech. Speak slowly and clearly, pausing for em[...] and to allow your audience to absorb what you're saying.

8. **Use storytelling techniques**: People remember stories better than facts and figures. Use anecdotes, ex[...] and metaphors to make your message more relatable and engaging.

9. **Anticipate questions**: Think about potential questions your audience may have and prepare responses.[...] help you feel more confident and in control.

10. **Join a public speaking group**: Consider joining a group like Toastmasters, where you can practice spea[...] supportive environment and receive feedback from others.

11. **Record yourself**: Record yourself speaking and watch the video to identify areas for improvement, such[...]

16.85 tokens/second

# Biggest limitation right now:
# No GPU support. Yet.

# My recommendation

While this is great for small tasks (classification, image tagging, etc.), any LLM work should be delegated to GPU based infrastructure.

# Your LLM options for now

Use a SaaS service: OpenAI, Claude API, etc. Or use HF Inference endpoints to deploy the model of your choice

And query that API endpoint from your app.

Or "Build your own", at your own risk!