# An Inaugural AI Journey with Langchain and



## The Monthly dev: April 30 2024

# About Michael Dawson

Node.js lead for Red Hat and IBM

Active Node.js community member

    Node.js Collaborator, Node.js Technical Steering Committee,

    Active in a number of Working group(s)

Active OpenJS Foundation member

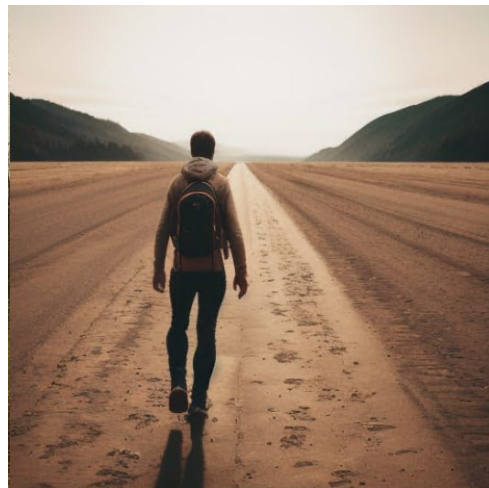    Voting Cross Project Council Member

    Community Director 2020-2022

Twitter: @mhdawson1

GitHub: @mhdawson

Linkedin: https://www.linkedin.com/in/michael-dawson-6051282

- Why Node.js ?
- My journey
  - Starting as a Newbie
  - Running a model locally
  - Leveraging a GPU
  - Retrieval Augmented Generation
  - Working with different model serving options

# Why Node.js ?

- Python often seen as runtime for AI
- But !!!
  - Not all applications will move to Python
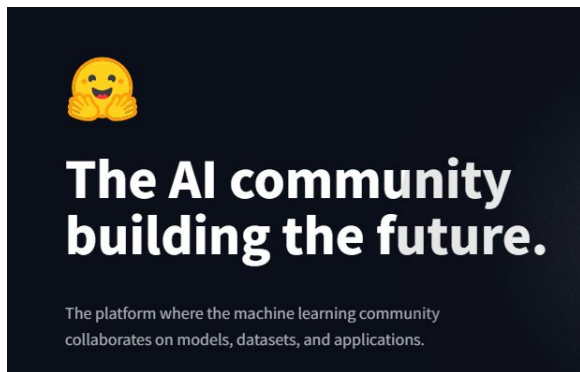  - Emerging AI client libraries often support TypeScript/JavaScript

- Until recently
  - https calls to bespoke service
- But, libraries are now emerging
  - langchain.com
  - llamaindex.ai
  -
  - .. add growing list here

- Need a model or remote API service
  - Need to be cautious with proprietary info
  - Choose to start by running locally

# Where do I get a model?

- [HuggingFace](#)

# How to I load a model?

- [llama.ccp](llama.ccp)
- [ollama](ollama)
- [Hugging Face transformers](Hugging Face transformers)
- …

- [node-llama-cpp](node-llama-cpp)

# What's this under the covers ?

- ## node-addon-api



**The state of the Node.js core: The Monthly Dev #29**

[Building Native addons like its 2023](#)

node-addon-api - https://github.com/nodejs/node-addon-api

# What's this under the covers ?

```cpp
if (info.Length() > 1 && info[1].IsObject()) {
    Napi::Object options = info[1].As<Napi::Object>();

    if (options.Has("gpuLayers")) {
        model_params.n_gpu_layers = options.Get("gpuLayers").As<Napi::Numb

    }

    if (options.Has("vocabOnly")) {
        model_params.vocab_only = options.Get("vocabOnly").As<Napi::Boolea

    }
```

# Where's the code ?



https://github.com/mhdawson/ai-experimentation

# Loading the model in Langchain.js

```javascript
//////////////////////////////
// GET THE MODEL
const __dirname = path.dirname(fileURLToPath(import.meta.url));
const modelPath = path.join(__dirname,
                            "models",
                            "mistral-7b-instruct-v0.1.Q5_K_M.gguf")
const { LlamaCpp } = await import("@langchain/community/llms/llama_cpp");
const model = await new LlamaCpp({ modelPath: modelPath });
```

# Asking my first question - create a chain

```
///////////////////////////////
// CREATE CHAIN
const prompt =
  ChatPromptTemplate.fromTemplate(`Answer the following question if you don't
know the answer say so:
Question: {input}`);
const chain = prompt.pipe(model);
```

# Asking my first question - ask a question

```javascript
//////////////////////////////////
// ASK QUESTION
console.log(new Date());
let result = await chain.invoke({
  input: "Should I use npm to start a node.js application",
});
console.log(result);
console.log(new Date());
```

# Asking my first question - 25 seconds later….

2024-03-11T22:08:23.372Z
Assistant: Yes, you should use npm to start a Node.js application. NPM (Node Package Manager) is the default package manager for Node.js and it provides a centralized repository of packages that can be used in your applications. It also allows you to manage dependencies between packages and automate tasks such as testing and deployment. If you are new to Node.js, I would recommend using npm to get started with your application development.
2024-03-11T22:08:45.774Z

Not the answer we want people to get based on the

nodejs-reference-architecture

https://github.com/nodeshift/nodejs-reference-architecture

**JSDrops Growing Success Across Organizations** -
https://www.youtube.com/watch?v=GncwXJBwcgQ

- Good news, node-llama-cpp supports GPUs
  - enabled by default for MacOS (non intel)
  - easy to enable for Windows

- [install the CUDA toolkit](install the CUDA toolkit) (version 12.x or higher).

- install the C/C++ compiler for your platform, including support for CMake and CMake.js.

- npx --no node-llama-cpp download --cuda

# Turning on the GPU

## 25 → 3 Seconds



NVIDIA 4060Ti 16G

- Often want to add additional knowledge
  - Building/Training a model is a lot of work
  - Just want to add a bit of specific info

# Fast but still wrong answer!

- Prompt Engineering
  - Prompt includes
    - question
    - context
  - For example in chatbot, context can include chat history

# Retrieval Augmented Generation (RAG)

- Ingest documents
- Extract relevant chunks
- Add chunks to prompt context

Note! - supported context is limited, for example 2k

# Load the documents

nodejs-reference-architecture

```javascript
const docLoader = new DirectoryLoader(
  "./SOURCE_DOCUMENTS",
  {
    ".md": (path) => new TextLoader(path),
  }
);
const docs = await docLoader.load();
```

# Split the documents

```
const splitter = await new MarkdownTextSplitter({
  chunkSize: 500,
  chunkOverlap: 50
});
const splitDocs = await splitter.splitDocuments(docs);
```

# Store the chunks in a database

```javascript
const vectorStore = await MemoryVectorStore.fromDocuments(
  splitDocs,
  new HuggingFaceTransformersEmbeddings()
);
const retriever = await vectorStore.asRetriever();
console.log("Augmenting data loaded - " + new Date());
```

```javascript
retriever.getRelevantDocuments("Should I use npm to start a node.js
application?");
```

# Create Chain

```
const prompt =
  ChatPromptTemplate.fromTemplate(`Answer the following question based only on the provided context, if you
don't know the answer say so:

<context>
{context}
</context>

Question: {input}`);

const documentChain = await createStuffDocumentsChain({
  llm: model,
  prompt,
});

const retrievalChain = await createRetrievalChain({
  combineDocsChain: documentChain,
  retriever,
});
```

# Ask Question

```javascript
/////////////////////////////////
// ASK QUESTIONS

console.log(new Date());
let result = await retrievalChain.invoke({
  input: "Should I use npm to start a node.js application",
});
console.log(result);
console.log(new Date());
```

# A Better Answer

'Assistant: It is generally not necessary to use `npm` to start a Node.js application. If you avoid using it in the container, you will not be exposed to any security vulnerabilities that might exist in that component or its dependencies. However, it is important to build security into your software development process when developing Node.js modules and applications. This includes managing dependencies, managing access and content of public and private data stores such as npm and github, writing defensive code, limiting required execution privileges, supporting logging and monitoring, and externalizing secrets.'

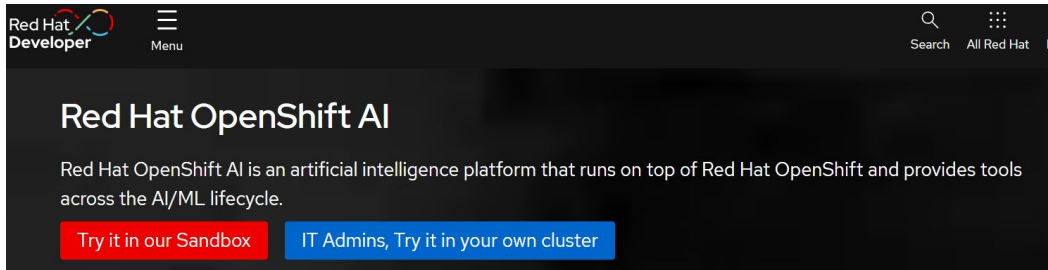# Switching to other model serving back ends

```
async function getModel(type, temperature) {
  console.log("Loading model - " + new Date());

  let model;
  if (type === 'llama-cpp') {
    const __dirname = path.dirname(fileURLToPath(import.meta.url));
    const modelPath = path.join(__dirname, "models", "mistral-7b-instruct-v0.1.Q5_K_M.gguf")
    const { LlamaCpp } = await import("@langchain/community/llms/llama_cpp");
    model = await new LlamaCpp({ modelPath: modelPath,
                                batchSize: 1024,
                                temperature: temperature,
                                gpuLayers: 64 });
  } else if (type === 'openAI') {
    const { ChatOpenAI } = await import("@langchain/openai");
    const key = await import('../key.json', { with: { type: 'json' } });
    model = new ChatOpenAI({
      temperature: temperature,
      openAIApiKey: key.default.apiKey
    });
  } else if (type === 'Openshift.ai') {
```
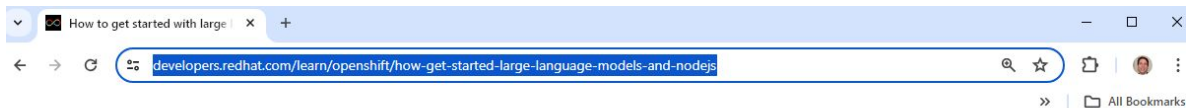
# Switching to other model serving back ends

```
} else if (type === 'openAI') {
  const { ChatOpenAI } = await import("@langchain/openai");
  const key = await import('../key.json', { with: { type: 'json' } });
  model = new ChatOpenAI({
    temperature: temperature,
    openAIApiKey: key.default.apiKey
  });
} else if (type === 'Openshift.ai') {
  const { ChatOpenAI } = await import("@langchain/openai");
  model = new ChatOpenAI(
    { temperature: temperature,
      openAIApiKey: 'EMPTY',
      modelName: 'mistralai/Mistral-7B-Instruct-v0.2' },
    { baseURL: 'http://vllm.llm-hosting.svc.cluster.local:8000/v1' }
  );
```

https://developers.redhat.com/products/red-hat-openshift-ai/overview

# To Dive Deeper

https://developers.redhat.com/learn/openshift/how-get-started-large-language-models-and-nodejs

# Some Key Takeaways

- It's easy to run a local model
- Strengths of Node.js still apply
- Libraries to improve DevX
  - TypeScript/JavaScript often supported

# Copyright and Trademarks