# AI Meets node.js

## Unlocking the Power of Large Language Models in JavaScript

POWERUp 2025

# 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
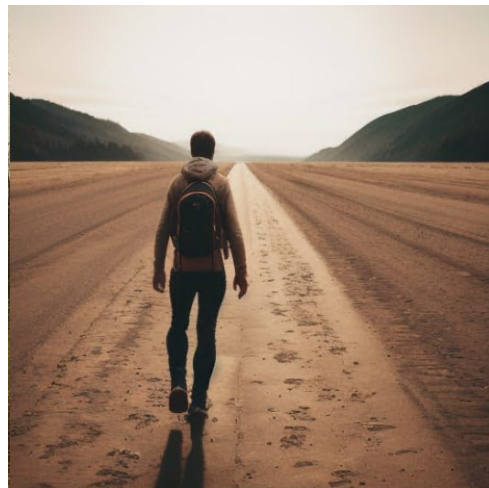
Twitter: @mhdawson1

GitHub: @mhdawson

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

- ## Why Node.js ?

- ## Our journey
  - Running a model locally
  - Leveraging a GPU
  - Some additional fundamentals, including Retrieval Augmented Generation
  - Working with different model serving options
  - Instrumenting the LLM components
  - A sample application - Parasol
  - A quick start with Podman AI Lab

# 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

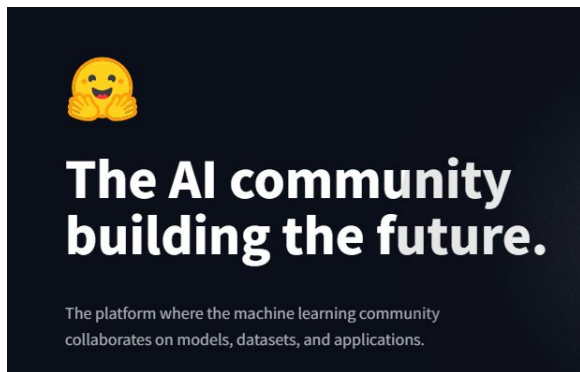- In the beginning
  - https calls to bespoke service
- But now, many libraries are now available
  - langchain.com/LangGraph
  - llamaindex.ai
  - ollama
  - Bee Agent
  - Llama Stack

- Need a model or remote API service
  - Need to be cautious with proprietary info
  - Choose to start by running locally
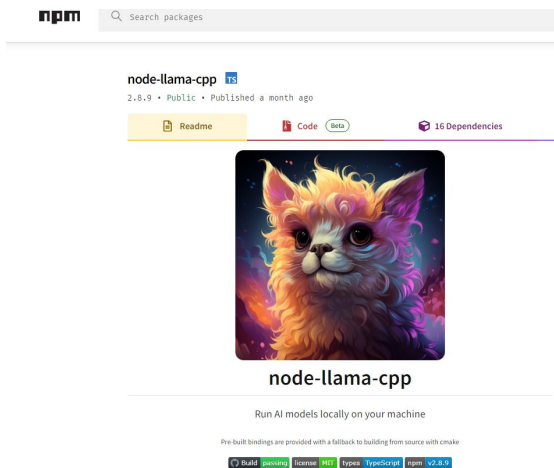
# Where do I get a model?

- ## HuggingFace



The AI community building the future.

The platform where the machine learning community collaborates on models, datasets, and applications.
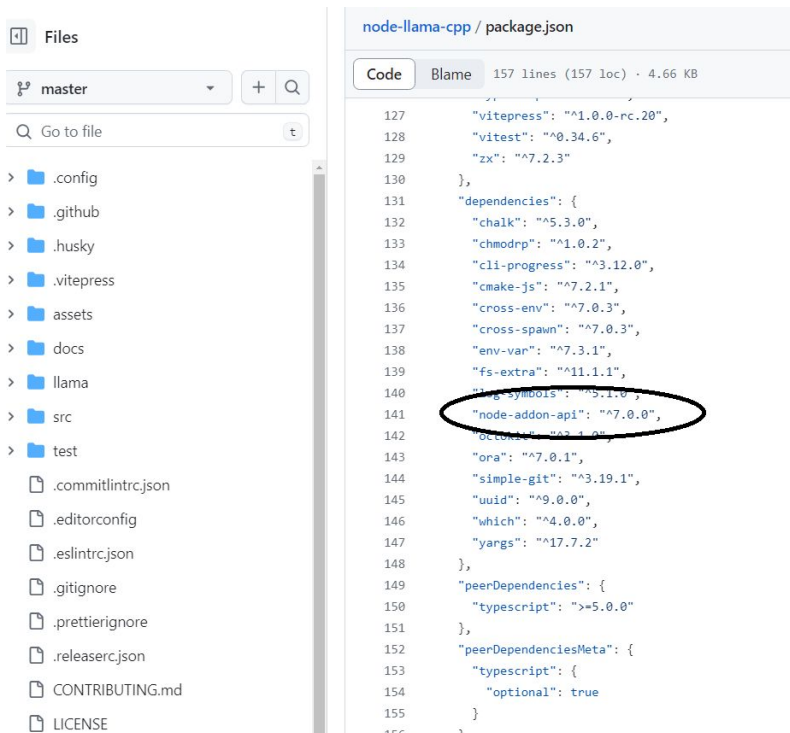
Be careful of which model you use!

# How do I load a model?

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

- [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::Number>().Int32Value
    }

    if (options.Has("vocabOnly")) {
        model_params.vocab_only = options.Get("vocabOnly").As<Napi::Boolean>().Value();
    }
```

# 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

```
//////////////////////////////////
// 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

# Hmm, 25 seconds is a bit long

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

# Turning on the GPU windows/NVIDIA

- [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

## 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

# Some Fundamentals - Overview

- Additional Context
- Message History
- Function Calling

# Retrieval Augmented Generation (RAG)

- Indexing
  - Load and Store Documents

  - Happens "Offline"

- Retrieval and Generation
  - Extract relevant chunks

  - Add chunks to prompt context

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

# Load the documents

nodejs-reference-architecture

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

# Split the documents

```javascript
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()
);
```

# Create Chain and Retrieve Docs

```javascript
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,
});
```

```javascript
const retriever = await vectorStore.asRetriever();


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.'

# Message History

- LLMs are Stateless
  - No context of previous questions

Example:

- Human: My name is Luke
- AI: Hello, Luke.
- Human: What's my name?
- AI: Huh? You never told me.

# Message History

```javascript
const prompt = ChatPromptTemplate.fromMessages([
  [ 'system', 'You are an assistant….'],
  new MessagesPlaceholder('history'),
  [ 'human', '{input}' ]
]);

const sessions = {};



const chainWithHistory = new RunnableWithMessageHistory({
  runnable: chain,
  getMessageHistory: (sessionId) => {...},
  inputMessagesKey: 'input',
  historyMessagesKey: 'history',
});
```

```javascript
const response1 = await chainWithHistory.invoke(
  { input: 'My name is Luke' },
  {
    configurable: {
      sessionId: 'funtimes'
    }
  }
);

console.log('Response 1: ', response1);



const response2 = await chainWithHistory.invoke(
  { input: 'What is my name' },
  {
    configurable: {
      sessionId: 'funtimes'
    }
  }
);

console.log('Response 2: ', response2);

console.log(sessions);
```

https://github.com/lholmquist/ascend_ai_talk/blob/main/message_history/history.js

# Message History

Response 1:   Hello, Luke! How can I assist you today?


Response 2:   Your name is Luke. Is there something specific you would like help with?

```
{
  funtimes: InMemoryChatMessageHistory {
    messages: [
      HumanMessage {
        "content": "My name is Luke",
      },
      AIMessage {
        "content": " Hello, Luke! How can I assist you
today?",
      },
      HumanMessage {
        "content": "What is my name",
      },
      AIMessage {
        "content": " Your name is Luke. Is there something
specific you would like help with?"
      }
    ]
  }
}
```

# Agents and Function Calling

- Respond by "calling a tool"
- LLM doesn't run the tool
  - generates the arguments

- Not Universal
  - Supported by many popular LLM providers



```
AIMessage(
  tool_calls=[{
    name: "get_weather"
    args: {
      location: "Hawaii",
    },
    id: "call_abc123",
    type: "tool_call"
  }]
)
```

https://js.langchain.com/docs/how_to/tool_calling/

# Agents and Function Calling

```javascript
const calculatorTool = tool(
  async ({ operation, number1, number2 }) => {
    if (operation === "add") {
      return `${number1 + number2}`;
    } else if (operation === "subtract") {
      return `${number1 - number2}`;
    } else if (operation === "multiply") {
      return `${number1 * number2}`;
    } else if (operation === "divide") {
      return `${number1 / number2}`;
    } else {
      throw new Error("Invalid operation.");
    }
  },
  {
    name: "calculator",
    description: "Can perform mathematical operations.",
    schema: calculatorSchema,
  }
);

const llmWithTools = model.bindTools([calculatorTool]);
```

https://github.com/lholmquist/ascend_ai_talk/blob/main/tools/tools.js

```javascript
const messages = [new HumanMessage("What is 3 * 12?")];

const aiMessage = await llmWithTools.invoke(messages);

AIMessage {
  "content": "",
  "additional_kwargs": {
    "tool_calls": [
      { "id": "call_kfmpr9ux",
        "type": "function" },
    ]
  },
  "tool_calls": [
    {
      "name": "calculator",
      "args": {
        "number1": 3,
        "number2": 12,
        "operation": "multiply"
      },
      "type": "tool_call",
      "id": "call_kfmpr9ux"
    },
  ],
}
```

# Agents and Function Calling

```
for (const toolCall of aiMessage.tool_calls) {
  const toolMessage = await calculatorTool.invoke(toolCall);

  console.log(toolMessage);

  messages.push(toolMessage);

}
```

```
Value of toolCall:
{
  name: 'calculator',
  args: {
    number1: 3,
    number2: 12,
    operation: 'multiply'
  },
  type: 'tool_call',
  id: 'call_kfmpr9ux'
}
```

```
    ToolMessage {
        "content": "36",
        "name": "calculator",
        "additional_kwargs": {},
        "response_metadata": {},
        "tool_call_id": "call_kfmpr9ux"
    },
```

# Agents and Function Calling

```
[
  HumanMessage {
    "content": "What is 3 * 12?"
  },
  AIMessage {
    "id": "chatcmpl-352",
    "content": "",
    "additional_kwargs": {
      "tool_calls": [
        {
          "id": "call_kfmpr9ux",
          "type": "function"
        },
      ]
    },
    "tool_calls": [
      {
        "name": "calculator",
        "args": {
          "number1": 3,
          "number2": 12,
          "operation": "multiply"
        },
        "type": "tool_call",
        "id": "call_kfmpr9ux"
      },
    ],
  },
  ToolMessage {
    "content": "36",
    "name": "calculator",
    "tool_call_id": "call_kfmpr9ux"
  },
]
```

```javascript
const result = await llmWithTools.invoke(messages);
console.log(result);
```

```
AIMessage {
  "id": "chatcmpl-497",
  "content": "36 is the result of multiplying 3 by 12"
  },
  "tool_calls": []
}
```



[LangGraph.js](LangGraph.js)

# Switching to other model serving back ends

```javascript
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
    });
```

https://github.com/mhdawson/ai-experimentation/blob/main/lesson-5/langchainjs-ollama.mjs

# Switching to other model serving back ends

```
  } else if (type === 'Openshift.ai') {
  ////////////////////////////////
  // Connect to OpenShift.ai endpoint
  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' }
  );
} else if (type === 'ollama') {
  ////////////////////////////////
  // Connect to ollama endpoint
  const { Ollama } = await import("@langchain/community/llms/ollama");
  model = new Ollama({
     baseUrl: "http://10.1.1.39:11434", // Default value
     model: "mistral", // Default value
  });
};
```



Red Hat Developer

Red Hat OpenShift AI

Red Hat OpenShift AI is an artificial intelligence platform that runs on top of Red Hat OpenShift and provides tools across the AI/ML lifecycle.

Try it in our Sandbox    IT Admins, Try it in your own cluster

# Instrumenting the LLM components

- ## OpenTelemetry
  - ### Trending to become common standard
- ## Red Hat build of OpenTelemetry

- ## LLM instrumentation following the trend
  - ### openllmetry
  - ### langtrace



**What do we instrument?**

OpenLLMetry-JS can instrument everything that OpenTelemetry already instruments - so things like your DB, API calls, and more. On top of that, we built a set of custom extensions that instrument things like your calls to OpenAI or Anthropic, or your Vector DB like Pinecone, Chroma, or Weaviate.

**LLM Providers**

- ✅ OpenAI
- ✅ Azure OpenAI
- ✅ Anthropic
- ✅ Cohere
- ⏳ Replicate
- ⏳ HuggingFace
- ✅ Vertex AI (GCP)
- ✅ Bedrock (AWS)

**Vector DBs**

- ✅ Pinecone
- ✅ Chroma
- ⏳ Weaviate
- ⏳ Milvus

**Frameworks**

- ✅ LangChain
- ✅ LlamaIndex

https://github.com/traceloop/openllmetry-js#-what-do-we-instrument

# Instrumenting the LLM components

## Auto Instrumentation

```
import * as traceloop from "@traceloop/node-server-sdk";
import { trace, context } from "@opentelemetry/api";
import * as LlamaIndex from "llamaindex";
import { OTLPTraceExporter } from "@opentelemetry/exporter-trace-otlp-grpc";
import { ConsoleSpanExporter } from "@opentelemetry/sdk-trace-node";

traceloop.initialize({ exporter: new OTLPTraceExporter(),          ⟵          Often Optional
                       disableBatch: true,
                       instrumentModules: {
                         llamaIndex: LlamaIndex,
                       }
                     });

const tracer = trace.getTracer();

import {askQuestions} from "./llamaindex-function-ollama.mjs"
tracer.startActiveSpan('Asking questions', async (span) => {
  await askQuestions();
  span.end();
});
```

# Instrumenting the LLM components

# Instrumenting the LLM components

# IBM i as Client

- favorite-color.mjs - https://github.com/mhdawson/ai-tool-experimentation/blob/main/ollama/favorite-color.mjs
- Tests function calling
  - favoriteColorTool
  - favoriteHockeyTeamTool
- Questions

  'What is my favorite color?',
  'My city is Ottawa',
  'My country is Canada',
  'I moved to Montreal. What is my favorite color now?',
  'My city is Montreal and my country is Canada',
  'My city is Ottawa and my country is Canada, what is my favorite color?',
  'What is my favorite hockey team ?',
  'My city is Montreal and my country is Canada',

# IBM i as Client

Iteration 0 -----------------------------------------------------------
QUESTION: What is my favorite color?
  CALLED:favoriteColorTool
  RESPONSE:To determine your favorite color, I need to know where you're from. Can you tell me which city and country you are in?
QUESTION: My city is Ottawa
  CALLED:favoriteColorTool
  RESPONSE:Based on our tool's analysis, it seems like your favorite color is black. Is that correct?
QUESTION: My country is Canada
  CALLED:favoriteColorTool
  RESPONSE:Since you provided your country as Canada, but I asked for a city earlier and you said Ottawa, I'll try again. Since you already mentioned Ottawa earlier, let's just go with that.

So, considering your favorite color tool analysis for Ottawa, Canada, it seems like your favorite color is black. Is that correct?
QUESTION: I moved to Montreal. What is my favorite color now?
  CALLED:favoriteColorTool
  RESPONSE:Now that you've moved to Montreal, it seems like your favorite color has changed! According to our tool's analysis, your new favorite color is actually red. Is that a change you're happy with?
QUESTION: My city is Montreal and my country is Canada
  CALLED:favoriteHockeyTeamTool
  RESPONSE:Now that you've provided both your city and country, I can see that our tool has actually changed to a different function. It seems like we're not looking at favorite colors anymore, but rather favorite hockey teams!

…

# GPUs on IBM i

- Today:
  - on-chip acceleration with MMA (mostly optimized for

    deep learning and ML inferencing but also has cost/perf benefits over GPUs for small LLMs)
  - available on Power 10+
- Soon:
  - IBM Telum chip on [Spyre](#) cards for genAI acceleration
  - available with Power 11+

# A sample application - Parasol

7. Definitions

   7.1. "Insured vehicle" refers to the automobile listed on the declarations page.

   7.2. "Accident" means a sudden, unexpected event resulting in damage or injury.

8. Additional Provisions

   8.1. Coverage extends to other drivers listed on the policy.

   8.2. Rental car coverage may be included if specified on the declarations page.

   8.3. Roadside assistance may be available if specified on the declarations page.

9. Contact Information

For claims or inquiries:

Phone: 800-CAR-SAFE

Email: claims@parasol.com

# To learn more about Parasol

- [Improving Chatbot result with Retrieval Augmented Generation (RAG) and Node.js | Red Hat Developer](#)
- [Experimenting with Email generation and summarization with Node.js and Large Language Models | Red Hat Developer](#)
- [Chatbot, I Choose You….. to call a function | Red Hat Developer](#)
- [Exploring an insurance use case with AI and Node.js | Red Hat Developer](#)

# Podman Desktop



Containers and Kubernetes for application developers

Podman Desktop is an open source graphical tool enabling you to seamlessly work with containers and Kubernetes from your local environment.

Download Now

For **Windows** *(browser-detected)*

Other downloads



Podman AI Lab

# Podman AI lab

# Node.js RAG Recipe

# To Dive Deeper

[How to get started with large language models and Node.js | Red Hat Developer](#)
[Diving Deeper with large language models and Node.js | Red Hat Developer](#)
[Essential AI tutorials for Node.js Developers](#)





[AI & Node.js | Red Hat Developer](#)
[A Developer's Guide to the Node.js Reference Architecture](#)

# Copyright and Trademarks