



Artificial Intelligence Using Swift

by Mark Watson



Artificial Intelligence Using Swift

CoreML, NLP, Deep Learning, Semantic Web and Linked Data, Knowledge Graphs, Knowledge Representation

Mark Watson

This book is available at <https://leanpub.com/SwiftAI>

This version was published on 2025-06-26



This is a [Leanpub](#) book. Leanpub empowers authors and publishers with the Lean Publishing process. [Lean Publishing](#) is the act of publishing an in-progress ebook using lightweight tools and many iterations to get reader feedback, pivot until you have the right book and build traction once you do.

© 2020 - 2025 Mark Watson

Contents

Cover Material, Copyright, and License	1
Preface	2
Notes on the new June 2025 Book Edition	2
Requests from the Author	2
Notes on the Second Edition	2
Book Structure	3
Requests from the Author	3
Parts of this Book are Specific for macOS and iOS, with Some Support for Linux	4
Code for this Book	4
Author's Background	5
Cover Art	5
CoreML Libraries Used in this Book	5
Swift 3rd Party Libraries	6
Acknowledgements	6
Part 1: Introduction and Short Examples	7
Setting Up Swift for Command Line Development	8
Installing Swift Packages	8
Creating Swift Packages	9
Accessing Libraries that You Write in Other Projects	9
Background Information for Writing Swift Command Line Utilities	13
Using Shell Processes	13
FileIO Examples	16
Swift REPL	18
Web Scraping	20
Running in the Swift REPL	24
Part 2: Large Language Models	26
Using the OpenAI LLM APIs	27
Core Architecture	30
Key Features	30

CONTENTS

Technical Implementation Details	31
Running Tests	32
Using APIs for Anthropic Claude LLMs	34
Running the examples	34
Using Groq APIs to Open Weight LLM Models	37
Implementation of a Client Library for the Groq APIs	37
Running the Tests	41
Using the xAI Grok LLM	43
Implementation of a Grok API Client Library	43
Using Ollama to Run Local LLMs	48
Running the Ollama Service	48
Ollama Wrap Up	50
Using Apple's MLX Framework to Run Local LLMs	51
MLX Framework History	51
MLX Resources on GitHub	51
Example Application for MLX Swift Examples Repository	52
Analysis of Swift and SwiftUI Code in the LLMEval Application	55
Part 3: Apple's CoreML and NLP Libraries	62
Deep Learning Introduction	63
Simple Multi-layer Perceptron Neural Networks	64
Deep Learning	66
Natural Language Processing Using Apple's Natural Language Framework	67
Using Apple's NaturalLanguage Swift Library	67
A simple Wrapper Library for Apple's NLP Models	67
Documents Question Answering Using OpenAI GPT4 APIs and a Local Embeddings Vector Database	71
Extending the String Class	71
Implementing a Local Vector Database for Document Embeddings	72
Create Local Embeddings Vectors From Local Text Files With OpenAI GPT APIs	74
Using Local Embeddings Vector Database With OpenAI GPT APIs	76
Wrap Up for Using Local Embeddings Vector Database to Enhance the Use of GPT3 APIs With Local Documents	79
Part 4: Knowledge Representation and Data Acquisition	80
Linked Data and the Semantic Web	81
Understanding the Resource Description Framework (RDF)	84

CONTENTS

Frequently Used Resource Namespaces	84
Understanding the SPARQL Query Language	86
Semantic Web and Linked Data Wrap Up	86
Example Application: iOS and macOS Versions of my KnowledgeBookNavigator	87
Screen Shots of macOS Application	87
Application Code Listings	90
Part 5: Apple Intelligence	108
Developers Can Now Weave Apple Intelligence Directly Into Their Apps	108
Key Advantages for Developers:	108
Using Apple Intelligence's Default System Model To Build a Chat Command Line Tool	109
Using Apple Intelligence's Default System Model To Build a Coding Assistant Command Line Tool	113
Book Wrap Up	117

Cover Material, Copyright, and License

Copyright 2022-2025 Mark Watson. All rights reserved. This book may be shared using the Creative Commons “share and share alike, no modifications, no commercial reuse” license.

This eBook will be updated occasionally so please periodically check the [leanpub.com web page for this book¹](#) for updates.

The first edition was released spring of 2022. The second edition was released December 2024. The third edition was released June 2025.

If you would like to support my work please consider purchasing my books on [Leanpub²](#) and star my git repositories that you find useful on [GitHub³](#). You can also interact with me on social media on [Mastodon⁴](#) and [Twitter⁵](#).

¹<https://leanpub.com/SwiftAI>

²<https://leanpub.com/u/markwatson>

³<https://github.com/mark-watson?tab=repositories&q=&type=public>

⁴https://mastodon.social/@mark_watson

⁵https://twitter.com/mark_l_watson

Preface

Why use Swift for hacking AI? Common Lisp has been my go-to language for artificial intelligence development and research since 1982. The transition to using Swift was a transition motivated by practical aspects of Swift and the Swift ecosystem.

Notes on the new June 2025 Book Edition

With the release of new Apple Intelligence tooling available for macOS26, iPadOS26, and iOS26 I have added a new Part 5 at the end of this book covering several practical AI use cases with new examples.

Requests from the Author

This book will always be available to read free online at <https://leanpub.com/SwiftAI/read>⁶.

That said, I appreciate it when readers purchase my books because the income enables me to spend more time writing.

Hire the Author as a Consultant

I am available for short consulting projects. Please see <https://markwatson.com>⁷.

You can also interact with me on social media on [Mastodon](#)⁸ and [Twitter](#)⁹.

Notes on the Second Edition

The second edition of this book deletes some of the old material and adds two new themes:

- A new Part II of the book that covers Large Language Models (LLMs). We will use both commercial LLM APIs and running local LLMs using Ollama and Apple's MLX framework.
- Several examples from the first edition are augmented using LLMs.
- As much as possible, I support some of the book examples as Swift Playgrounds, usable on iPads and Macs.

⁶<https://leanpub.com/SwiftAI/read>

⁷<https://markwatson.com>

⁸[@mark_watson](https://mastodon.social/@mark_watson)

⁹https://twitter.com/mark_l_watson

Book Structure

This book starts out slowly in Part I with simple examples which I wrote showing how to access the Swift library packages on GitHub, tips on writing Swift command line apps.

Part II will show you how to effectively integrate LLMs into your own applications.

Part III starts with a simple example using web scraping and commercial web search APIs. We then work through examples integrating web search with LLMs and then show how we can modify web scraping applications to specifically process topics and have better control of outputting structured data.

We then proceed to using Apple's CoreML for Natural Language Processing (NLP), training and using your own CoreML models, using OpenAI's GPT-4 APIs, and finally several semantic web/linked data examples. The book ends with the example macOS application Knowledge Graph Navigator. It is not my intention to cover in detail the use of SwiftUI for building iOS/iPadOS/macOS applications but I thought my readers might enjoy seeing several of the techniques covered in the book integrated into an example app.

I have used Common Lisp for AI research projects and for AI product development and delivery since 1982. There is something special about using a language for almost forty years. I now find Swift a compelling choice for several reasons:

- Flexible language with many features I rely on like supporting closures and an interactive functional programming style.
- Built-in support for deep learning neural network models for natural language processing, predictive models, etc.
- First class support for iOS and macOS development.
- Good support for server side applications hosted on Linux.

Swift is a programmer-efficient language: code is concise and easy to read, and high quality libraries from Apple and third parties mean that often there is less code to write. I will share with you my Swift development work flow that combines interactive development of code in playgrounds, development of higher level libraries in text only or command line applications, and my general strategy for writing iOS and macOS applications after low level and intermediate code is written and debugged.

Requests from the Author

This book will always be available to read free online at <https://leanpub.com/SwiftAI/read>¹⁰.

That said, I appreciate it when readers purchase my books because the income enables me to spend more time writing.

¹⁰<https://leanpub.com/SwiftAI/read>

Hire the Author as a Consultant

I am available for short consulting projects. Please see <https://markwatson.com>¹¹.

Parts of this Book are Specific for macOS and iOS, with Some Support for Linux

Swift is a general purpose language that is well supported in macOS, iOS, and Linux, with some support in Windows. Here, we cover the use of Swift on macOS and iOS. Some of the examples in this book rely on libraries that are specifically available on macOS and iOS like CoreML and the NLP libraries. Several book examples also work on Linux, such as the examples using SQLite, the Microsoft Azure search APIs, web scraping, and semantic web/linked data.

Code for this Book

Because of the way the Swift Package Manager works, I organized all book examples that build libraries as separate GitHub repos so the libraries can be easily used in other book examples as well as your own software projects. The separate library GitHub repositories are:

- https://github.com/mark-watson/SparqlQuery_swift¹² - SPARQL Swift library for my Swift AI book.
- https://github.com/mark-watson/QuestionAnswering_BERT_swift¹³ - modification of Apple's question answering demo to use DBpedia.
- https://github.com/mark-watson/swift-coreml-wisconsin_data_create_model¹⁴ - create CoreML models from training data files of Wisconsin Cancer data.
- https://github.com/mark-watson/swift-coreml-wisconsin_data_predict_with_model¹⁵ - use the pretrained Wisconsin Cancer data model.
- https://github.com/mark-watson/ShellProcess_swift¹⁶ - library for spawning shell processes and capturing output to stdout.
- https://github.com/mark-watson/WebScraping_swift¹⁷ - library for scrapping web sites.
- https://github.com/mark-watson/OpenAI_swift¹⁸ - library for using OpenAI's GPT3 APIs.
- https://github.com/mark-watson/Nlp_swift¹⁹ - library that uses pretrained CoreML NLP models.

¹¹<https://markwatson.com>

¹²https://github.com/mark-watson/SparqlQuery_swift

¹³https://github.com/mark-watson/QuestionAnswering_BERT_swift

¹⁴https://github.com/mark-watson/swift-coreml-wisconsin_data_create_model

¹⁵https://github.com/mark-watson/swift-coreml-wisconsin_data_predict_with_model

¹⁶https://github.com/mark-watson/ShellProcess_swift

¹⁷https://github.com/mark-watson/WebScraping_swift

¹⁸https://github.com/mark-watson/OpenAI_swift

¹⁹https://github.com/mark-watson/Nlp_swift

- <https://github.com/mark-watson/KGN²⁰> - SwiftUI based application supporting macOS, iPadOS, and iOS. The macOS version is in Apple's app store.

I suggest cloning all of these GitHub repositories right now so you can have the example source code at hand while reading this book.

All of the code examples are licensed using the Apache 2 license. You are free to reuse the book example code in your own projects (open source, commercial), with attribution of my copyright and the Apache 2 license.

Except for the last SwiftUI example application, all sample programs are written as command line utilities. I considered using Swift playgrounds for some of the examples but decided that packaging as a combination of libraries and command line utilities would tend to make the example code more useful for your own projects.

<http://www.knowledgegraphnavigator.com/>

Author's Background

I have written 20+ books, mostly about artificial intelligence. I have over 50 US patents.

I write about technologies that I have used throughout my career: knowledge representation using semantic web and linked data, machine learning and deep learning, and natural language processing. I am grateful for the companies where I have worked (SAIC, Google, Capital One, Olive AI, Babylot, etc.) that have supported this work since 1982.

As an author, I hope that the material in this book entertains you and will be useful in your work.

Cover Art

The cover picture was taken by [WikiMedia Commons user Keta²¹](#) and is available for use under the Creative Commons License CC BY-SA 2.5.

CoreML Libraries Used in this Book

- CoreML general overview: <https://developer.apple.com/documentation/coreml>
- MLClassifier <https://developer.apple.com/documentation/createml/mlclassifier>
- MLTextClassifier <https://developer.apple.com/documentation/createml/mltextclassifier>
- NLModel <https://developer.apple.com/documentation/naturallanguage/nlmodel>
- Natural Language Framework <https://developer.apple.com/documentation/naturallanguage>
- MLCustomLayer <https://developer.apple.com/documentation/coreml/mlcustomlayer>

²⁰<https://github.com/mark-watson/KGN>

²¹<https://commons.wikimedia.org/wiki/User:Keta>

Swift 3rd Party Libraries

We use the following 3rd party libraries:

- <https://github.com/SwiftyJSON/SwiftyJSON>²²

Acknowledgements

I thank my wife Carol for editing this manuscript, finding typos, and suggesting improvements.

²²<https://github.com/SwiftyJSON/SwiftyJSON>

Part 1: Introduction and Short Examples

We begin with a sufficient introduction for Swift to understand the programming examples. After introducing the language we will look at a few short examples that provide code and techniques we use later in the book:

- Creating Swift projects
- Writing command line utilities
- Web scraping

Setting Up Swift for Command Line Development

Except for the last chapter in this book that uses Xcode for developing a complete macOS/iOS/iPadOS example application, I assume that you will work through the book examples using the command line and your favorite editor. If you want to use Xcode for the command line examples, you can open the Swift package file on the command line and open Xcode using, for example:

```
cd SparqlQuery_swift  
open Package.swift
```

You notice that most of the examples are command line apps or libraries with command line test programs and the **README.md** files in the example directories provide instructions for building and running on the command line.

You can also run Xcode and from the File Menu open an example **Package.swift** file. You can then use the Product / Test menu to run the test code for the example. You might need to use the View / Debug Area / Active Console menu to show the output area.

I assume that you are familiar with the Swift programming language and Xcode.

Swift is a general purpose language that is well supported in macOS and iOS, with good support for Linux, and with some support in Windows. For the purposes of this book, we are only considering the use of Swift on macOS and iOS. Most of the examples in this book rely on libraries that are specifically available on macOS and iOS like CoreML and the NLP libraries.

There are great free resources for the Swift language on the web, in other commercial books, and Apple's free Swift books. Here I provide just enough material on the Swift language for you to understand and work with the book examples. After working through this book's material you will be able to add machine learning, natural language processing, and knowledge representation to your applications. There will be parts of the Swift language that we don't need for the material here, and we won't cover.

Installing Swift Packages

We will use the [Swift Package Manager²³](#). You should pause reading now and install the Swift Package Manager if you have not already done so.

²³<https://swift.org/package-manager/>

I occasionally use <https://vapor.codes> web framework²⁴ library (although not in this book). We use this 3rd party library as an example for building a library locally from source code. Start by cloning the git repository <https://github.com/vapor/vapor>²⁵. Then:

```
git clone https://github.com/vapor/vapor.git  
cd vapor  
swift build
```

I don't usually install libraries locally from source code unless I am curious about the implementation and want to read through the source code. Later we will see how to reference Swift libraries hosted on GitHub in a project's `Package.swift` file.

Creating Swift Packages

We will cover using the Swift Package Manager to create new packages using the command line here. Later we will create projects using Apple's XCode IDE when we develop the example application Knowledge Graph Navigator.

You will want to use the [Swift Package Manager documentation](#)²⁶ for reference.

We will be generating executable projects and library (with a sample main program) projects. The commands for generating the stub for an executable application project are:

```
mkdir BingSearch  
cd BingSearch  
swift package init --type executable
```

and build the stub of a library with a demo main program:

```
mkdir SparqlQuery  
cd SparqlQuery  
swift package init --type library
```

Accessing Libraries that You Write in Other Projects

You can reference Swift libraries using the `Swift.package` file for each of your projects. We will look at parts of two `Swift.package` files here. The first is for my SPARQL query client library that we will develop in a later chapter. This library `SparqlQuery_swift` is used in both book examples Knowledge Graph Navigator (KGN) macOS/iOS/iPadOS example application as well as a text version `KnowledgeGraphNavigator_swift`.

²⁴<https://vapor.codes>

²⁵<https://github.com/vapor/vapor>

²⁶<https://github.com/apple/swift-package-manager/blob/main/Documentation/Usage.md>

```
1 import PackageDescription
2
3 let package = Package(
4     name: "SparqlQuery_swift",
5     products: [
6         .library(
7             name: "SparqlQuery_swift",
8             targets: ["SparqlQuery_swift"]),
9     ],
10    dependencies: [
11        .package(url: "https://github.com/SwiftyJSON/SwiftyJSON.git",
12            .branch("master")),
13    ],
14    targets: [
15        .target(
16            name: "SparqlQuery_swift",
17            dependencies: ["SwiftyJSON"]),
18        .testTarget(
19            name: "SparqlQuery_swiftTests",
20            dependencies: ["SparqlQuery_swift", "SwiftyJSON"]),
21    ]
22 )
```

This Swift package file is used to declare a Swift package named “SparqlQuery_swift”. The package contains one library target named “SparqlQuery_swift” and one test target named “SparqlQuery_swiftTests”. The library target depends on the “SwiftyJSON” package, which is specified as a dependency in the “dependencies” section of the package.

The “products” section defines the products that this package provides. In this case, the package provides a library product named “SparqlQuery_swift”. The library is built from the source code in the “SparqlQuery_swift” target.

The “dependencies” section lists the packages that this package depends on. In this case, it depends on the “SwiftyJSON” package, which is specified as a Git repository URL.

The “targets” section lists the targets that are part of the package. In this case, there are two targets: “SparqlQuery_swift” and “SparqlQuery_swiftTests”. The “SparqlQuery_swift” target depends on “SwiftyJSON”. The “SparqlQuery_swiftTests” target depends on both “SparqlQuery_swift” and “SwiftyJSON”.

The `Swift.package` file for text version `KnowledgeGraphNavigator_swift` is shown here:

```
1 import PackageDescription
2
3 let package = Package(
4     name: "KnowledgeGraphNavigator_swift",
5     platforms: [
6         .macOS(.v10_15),
7     ],
8     dependencies: [
9         .package(url: "https://github.com/SwiftyJSON/SwiftyJSON.git",
10            .branch("master")),
11         .package(url: "https://github.com/scinfinu/SwiftSoup.git", from: "1.7.4"),
12         .package(url: "git@github.com:mark-watson/SparqlQuery_swift.git",
13            .branch("main")),
14         .package(url: "git@github.com:mark-watson/Nlp_swift.git", .branch("main")),
15     ],
16     targets: [
17         // Targets are the basic building blocks of a package.
18         // A target can define a module or a test suite.
19         // Targets can depend on other targets in this package,
20         // and on products in packages this package depends on.
21         .target(
22             name: "KnowledgeGraphNavigator_swift",
23             dependencies: ["SparqlQuery_swift", "Nlp_swift",
24                           "SwiftyJSON", "SwiftSoup"]),
25     ]
26 )
```

This Swift package file is used to declare a Swift package named “KnowledgeGraphNavigator_swift”. The package contains one target named “KnowledgeGraphNavigator_swift”. The target depends on the “SparqlQuery_swift”, “Nlp_swift”, “SwiftyJSON”, and “SwiftSoup” packages, which are specified as dependencies in the “dependencies” section of the package.

The “platforms” section specifies the minimum platform version that the package supports. In this case, the package supports macOS version 10.15 and later.

The “dependencies” section lists the packages that this package depends on. In this case, it depends on four packages:

- SwiftyJSON: a Swift library for working with JSON data.
- SwiftSoup: a Swift library for parsing HTML and XML documents.
- SparqlQuery_swift: a Swift library for querying RDF data using the SPARQL query language.
- Nlp_swift: a Swift library for natural language processing.

The “targets” section lists the targets that are part of the package. In this case, there is one target named **KnowledgeGraphNavigator_swift**. The target depends on “**SparqlQuery_swift**, **Nlp_swift**, **SwiftyJSON**, and **SwiftSoup**.

Hopefully you have cloned the git repositories for each book example and understand how I have configured the examples for your use.

For the rest of this book, you can read chapters in any order. In some cases, earlier chapters will contain implementations of libraries used in later chapters.

Background Information for Writing Swift Command Line Utilities

This short chapter contains example code and utilities for writing command line programs, using external shell processes, and using the FileIO library.

Using Shell Processes

The library for using shell processes is one of my GitHub projects so you can include it in other projects using:

```
1 dependencies: [
2     .package(url: "git@github.com:mark-watson/ShellProcess_swift.git",
3             .branch("main")),
4 ] ,
```

You can clone this repository if you want to have the source code at hand:

```
1 git clone https://github.com/mark-watson/ShellProcess_swift.git
```

The following listing shows the library implementation. In line 5 we use the constructor `Process` from the Apple `Foundation` library to get a new process object that we set fields `executableURL` and `argList`. In lines 8 and 9 we create a new Unix style pipe to capture the output from the shell process we are starting and attach it to the process. After we run the task, we capture the output and return it as the value of function `run_in_shell`.

```
1 import Foundation
2
3 @available(OSX 10.13, *)
4 public func run_in_shell(commandPath: String, argList: [String] = []) -> String {
5     let task = Process()
6     task.executableURL = URL(fileURLWithPath: commandPath)
7     task.arguments = argList
8     let pipe = Pipe()
9     task.standardOutput = pipe
10    do {
```

```
11     try! task.run()
12     let data = pipe.fileHandleForReading.readDataToEndOfFile()
13     let output: String? = String(data: data, encoding: String.Encoding.utf8)
14     if let output = output {
15         if !output.isEmpty {
16             return output.trimmingCharacters(in: .whitespacesAndNewlines)
17         }
18     }
19 }
20 return ""
21 }
```

The function named `run_in_shell` takes two parameters: `commandPath` (a string representing the path to the executable command to be run) and `argList` (an array of strings representing the arguments to be passed to the command). The function returns a string that represents the output of the command.

Function `run_in_shell` first creates an instance of the `Process` class, which is used to run the command. It sets the `executableURL` property of the task instance to the `commandPath` value and sets the `arguments` property to the `argList` value. This function then creates a `Pipe` instance, which is used to capture the output of the command. It sets the `standardOutput` property of the task instance to the `Pipe` instance.

The function then runs the command using the `run()` method of the task instance. If the command runs successfully, the function reads the output of the command from the `Pipe` instance using the `readDataToEndOfFile()` method of the `fileHandleForReading` property. It then converts the output data to a string using the `String(data:encoding:)` initializer.

If the output string is not empty, this function trims leading and trailing whitespace and returns the resulting string. Otherwise, the function returns an empty string.

Overall, this function provides a simple way to run a shell command and capture its output in a Swift program.

As in most examples in this book we use the Swift testing framework to run the example code at the command line using `swift test`. Running `swift test` does an implicit `swift build`.

```
1 import XCTest
2 @testable import ShellProcess_swift
3
4 final class ShellProcessTests: XCTestCase {
5     func testExample() {
6         // This is an example of a functional test case.
7         // Use XCTAssert and related functions to verify your tests produce the
8         // correct results.
9         print("/** s1:")
10        let s1 = run_in_shell(commandPath: "/bin/ps", argList: ["a"])
11        print(s1)
12        let s2 = run_in_shell(commandPath: "/bin/ls", argList: ["."])
13        print("/** s2:")
14        print(s2)
15        let s3 = run_in_shell(commandPath: "/bin/sleep", argList: ["2"])
16        print("/** s3:")
17        print(s3)
18    }
19
20    static var allTests = [
21        ("testExample", testExample),
22    ]
23 }
24 }
```

This Swift unit test function is part of a test suite for the `ShellProcess_swift` package. The function is named `testExample` and is decorated with the `@testable` import statement to indicate that it tests an internal implementation detail of the `ShellProcess_swift` package.

The function uses the `run_in_shell` function to run three shell commands: `ps a`, `ls ..`, and `sleep 2`. It prints the output of each command to the console.

This test function is an example of a functional test case. It doesn't actually verify that the functions being tested produce the correct results. Instead, it's a simple way to visually inspect the output of the commands and ensure that they are working as expected.

The `allTests` variable is an array of tuples that map the test function names to the corresponding function references. This variable is used by the XCTest framework to discover and run the test functions.

The test output (with some text removed for brevity) is:

```
1 $ swift test
2 Test Suite 'All tests' started at 2021-08-06 16:36:21.447
3 ** s1:
4 PID TT STAT      TIME COMMAND
5 3898 s000 Ss    0:00.01 login -pf markw8
6 3899 s000 S+    0:00.18 -zsh
7 3999 s001 Ss    0:00.02 login -pfl markw8 /bin/bash -c exec -la zsh /bin/zsh
8 4000 s001 S+    0:00.38 -zsh
9 5760 s002 Ss    0:00.02 login -pfl markw8 /bin/bash -c exec -la zsh /bin/zsh
10 5761 s002 S    0:00.14 -zsh
11 8654 s002 S+   0:00.06 /Applications/Xcode.app/Contents/Developer/Toolchains/Xco\
12 deDefault.xctoolchain/usr/bin/swift-test
13 8665 s002 S    0:00.03 /Applications/Xcode.app/Contents/Developer/usr/bin/xctest\
14 /Users/markw_1/GIT_swift_book/ShellProcess_swift/.build/arm64-apple-macosx/debug/Sh\
15 ellProcess_swiftPackageTests.xctest
16 8666 s002 R    0:00.00 /bin/ps a
17 ** s2:
18 Package.swift
19 README.md
20 Sources
21 Tests
22 ** s3:
23
24 Test Suite 'All tests' passed at 2021-08-06 16:36:23.468.
25     Executed 1 test, with 0 failures (0 unexpected) in 2.019 (2.021) seconds
```

FileIO Examples

This file I/O example uses the `ShellProcess_swift` library we saw in the last section so if you were to create your own Swift project with the following code listing, you would have to add this dependency in the `Project.swift` file.

When writing command line Swift programs you will often need to do simple file IO so let's look at some examples here:

```

1 import Foundation
2 import ShellProcess_swift // my library
3
4 @available(OSX 10.13, *)
5 func test_files_demo() -> Void {
6     // In order to append to an existing file, you need to get a file handle
7     // and seek to the end of a file. The following will not work:
8     let s = "the dog chased the cat\n"
9     try! s.write(toFile: "out.txt", atomically: true,
10                 encoding: String.Encoding.ascii)
11    let s2 = "a second string\n"
12    try! s2.write(toFile: "out.txt", atomically: true,
13                 encoding: String.Encoding.ascii)
14    let aString = try! String(contentsOfFile: "out.txt")
15    print(aString)
16
17    // For simple use cases, simply appending strings, then writing
18    // the result atomically works fine:
19    var s3 = "the dog chased the cat\n"
20    s3 += "a second string\n"
21    try! s3.write(toFile: "out2.txt", atomically: true,
22                 encoding: String.Encoding.ascii)
23    let aString2 = try! String(contentsOfFile: "out2.txt")
24    print(aString2)
25
26    // list files in current directory:
27    let ls = run_in_shell(commandPath: "/bin/ls", argList: ["."])
28    print(ls)
29
30    // remove two temporary files:
31    let shellOutput = run_in_shell(commandPath: "/bin/rm",
32                                    argList: ["out.txt", "out2.txt"])
33    print(shellOutput)
34 }
35
36 if #available(OSX 10.13, *) {
37     test_files_demo()
38 }
```

The OS version checks in this Swift code use the `#available` conditional compilation block.

The **#available block** is used to conditionally compile code based on the availability of APIs or features in the operating system version. In this case, the code inside the `#available(OSX 10.13, *)` block will only be executed if the running operating system is macOS 10.13 or later.

If the running operating system version is earlier than 10.13, the code inside the `#available` block will be skipped and the program will exit without running the `test_files_demo()` function.

These operating system version checks are done to ensure that the program is only executed on operating systems that support the APIs and features used by the code. This helps to prevent runtime errors and crashes on older operating system versions that may not support the required features.

This function demonstrates how to write to and read from files using the `write(toFile:atomically:encoding:)` and `String(contentsOfFile:)` methods, how to list files in the current directory using the `ls` shell command, and how to remove files using the `rm` shell command.

I created a temporary Swift project with the previous code listing and a `Project.swift` file. I built and ran this example using the `swift` command line tool.

Unlike the example in the last section where we built a reusable library with a test program, here we have a standalone program contained in a single file so we will use `swift run` to build and run this example:

```
1 $ swift run
2 Fetching git@github.com:mark-watson/ShellProcess_swift.git from cache
3 Cloning git@github.com:mark-watson/ShellProcess_swift.git
4 Resolving git@github.com:mark-watson/ShellProcess_swift.git at main
5 [5/5] Build complete!
6 a second string
7
8 the dog chased the cat
9 a second string
10
11 Package.resolved
12 Package.swift
13 README.md
14 Sources
15 out.txt
16 out2.txt
```

Swift REPL

There is an example of using the Swift REPL at the end of the next chapter on web scraping. For reference, you can start a REPL with:

```
1 $ swift run --repl
2 Type :help for assistance.
3 1> import WebScraping_swift
4 2> webPageText(uri: "https://markwatson.com")
5 $R0: String = "Mark Watson: AI Practitioner and Polyglot Programmer"...
6 3> public func foo(s: String) -> String { return s }
7 4> foo(s: "cat")
8 $R1: String = "cat"
9 5>
```

You can import packages and interactively enter Swift expressions, including defining functions.

In the next chapter we will look at a longer example that scrapes web sites.

In the next chapter we will look at one more simple example, building a web scraping library, before getting to the machine learning and NLP part of the book.

Web Scraping

It is important to respect the property rights of web site owners and abide by their terms and conditions for use. This [Wikipedia article on Fair Use²⁷](#) provides a good overview of using copyright material.

The web scraping code we develop here uses the Swift library **SwiftSoup** that is loosely based on the BeautifulSoup libraries available in other programming languages.

For my work and research, I have been most interested in using web scraping to collect text data for natural language processing but other common applications include writing AI news collection and summarization assistants, trying to predict stock prices based on comments in social media which is what we did at Webmind Corporation in 2000 and 2001, etc.

I wrote a simple web scraping library that is available at [```
1 dependencies: \[
2 .package\(url: "git@github.com:mark-watson/WebScraping_swift.git",
3 .branch\("main"\)\),
4 \],
```](https://github.com/mark-watson/WebScraping_swift<sup>28</sup></a> that you can use in your projects by putting the following dependency in your <b>Project.swift</b> file:</p></div><div data-bbox=)

Here is the main implementation file for the library:

```
1 import Foundation
2 import SwiftSoup
3
4 public func webPageText(uri: String) -> String {
5 guard let myURL = URL(string: uri) else {
6 print("Error: \(uri) doesn't seem to be a valid URL")
7 fatalError("invalid URI")
8 }
9 let html = try! String(contentsOf: myURL, encoding: .ascii)
10 let doc: Document = try! SwiftSoup.parse(html)
11 let plain_text = try! doc.text()
12 return plain_text
13 }
14
```

---

<sup>27</sup>[https://en.wikipedia.org/wiki/Fair\\_use](https://en.wikipedia.org/wiki/Fair_use)

<sup>28</sup>[https://github.com/mark-watson/WebScraping\\_swift](https://github.com/mark-watson/WebScraping_swift)

```
15 func webPageHeadersHelper(uri: String, headerName: String) -> [String] {
16 var ret: [String] = []
17 guard let myURL = URL(string: uri) else {
18 print("Error: \(uri) doesn't seem to be a valid URL")
19 fatalError("invalid URI")
20 }
21 do {
22 let html = try String(contentsOf: myURL, encoding: .ascii)
23 let doc: Document = try SwiftSoup.parse(html)
24 let h1_headers = try doc.select(headerName)
25 for el in h1_headers {
26 let h1 = try el.text()
27 ret.append(h1)
28 }
29 } catch {
30 print("Error")
31 }
32 return ret
33 }
34
35 public func webPageH1Headers(uri: String) -> [String] {
36 return webPageHeadersHelper(uri: uri, headerName: "h1")
37 }
38
39 public func webPageH2Headers(uri: String) -> [String] {
40 return webPageHeadersHelper(uri: uri, headerName: "h2")
41 }
42
43 public func webPageAnchors(uri: String) -> [[String]] {
44 var ret: [[String]] = []
45 guard let myURL = URL(string: uri) else {
46 print("Error: \(uri) doesn't seem to be a valid URL")
47 fatalError("invalid URI")
48 }
49 do {
50 let html = try String(contentsOf: myURL, encoding: .ascii)
51 let doc: Document = try SwiftSoup.parse(html)
52 let anchors = try doc.select("a")
53 for a in anchors {
54 let text = try a.text()
55 let a_uri = try a.attr("href")
56 if a_uri.hasPrefix("#") {
57 ret.append([text, uri + a_uri])
58 }
59 }
60 } catch {
61 print("Error")
62 }
63 return ret
64 }
```

```
58 } else {
59 ret.append([text, a_uri])
60 }
61 }
62 } catch {
63 print("Error")
64 }
65 return ret
66 }
```

This Swift code defines several functions that can be used to scrape information from a web page located at a given URI.

The `webPageText` function takes a URI as input and returns the plain text content of the web page located at that URI. It first checks if the URI is valid and then reads the content of the web page using the `contentsOf` method of the `String` class. It then uses the `parse` method of the `SwiftSoup` library to parse the HTML content of the page and extract the plain text.

The `webPageH1Headers` and `webPageH2Headers` functions use the `webPageHeadersHelper` function to extract the H1 and H2 header texts respectively from the web page located at a given URI. The `webPageHeadersHelper` function uses the same technique as the `webPageText` function to read and parse the HTML content of the page. It then selects the headers using the specified `headerName` parameter and extracts the text of the headers.

The `webPageAnchors` function extracts all the anchor tags `<a>` from the web page located at a given URI, along with their corresponding text and URI. It also uses the `webPageHeadersHelper` function to read and parse the HTML content of the page, selects the anchor tags using the “`a`” selector, and extracts their `text` and `href` attributes.

Overall, these functions provide a simple way to scrape information from a web page and extract specific information such as plain text, header texts, and anchor tags.

I wrote these utility functions to get the plain text from a web site, HTML header text, and anchors. You can clone this library and extend it for other types of HTML elements you may need to process.

The test program shows how to call the APIs in the library:

```
1 import XCTest
2 import Foundation
3 import SwiftSoup
4
5 @testable import WebScraping_swift
6
7 final class WebScrapingTests: XCTestCase {
8 func testGetWebPage() {
9 let text = webPageText(uri: "https://markwatson.com")
```

```
10 print("\n\n\tTEXT FROM MARK's WEB SITE:\n\n", text)
11 }
12
13 func testToShowSwiftSoupExamples() {
14 let myURLString = "https://markwatson.com"
15 let h1_headers = webPageH1Headers(uri: myURLString)
16 print("\n\n++ h1_headers:", h1_headers)
17 let h2_headers = webPageH2Headers(uri: myURLString)
18 print("\n\n++ h2_headers:", h2_headers)
19 let anchors = webPageAnchors(uri: myURLString)
20 print("\n\n++ anchors:", anchors)
21 }
22
23 static var allTests = [("testGetWebPage", testGetWebPage),
24 ("testToShowSwiftSoupExamples",
25 testToShowSwiftSoupExamples)]
26 }
```

This Swift test program tests the functionality of the `WebScraping_swift` library. It defines two test functions: `testGetWebPage` and `testToShowSwiftSoupExamples`.

The `testGetWebPage` function uses the `webPageText` function to retrieve the plain text content of my website located at “<https://markwatson.com>”. It then prints the retrieved text to the console.

The `testToShowSwiftSoupExamples` function demonstrates the use of `webPageH1Headers`, `webPageH2Headers`, and `webPageAnchors` functions on the same website. It extracts and prints the H1 and H2 header texts and anchor tags of the same website.

The `allTests` variable is an array of tuples that map the test function names to the corresponding function references. This variable is used by the XCTest framework to discover and run the test functions.

Overall, this Swift test program demonstrates how to use the functions defined in the `WebScraping_swift` library to extract specific information from a web page.

Here we run the unit tests (with much of the output not shown for brevity):

```
1 $ swift test
2
3 TEXT FROM MARK's WEB SITE:
4
5 Mark Watson: AI Practitioner and Polyglot Programmer | Mark Watson Read my Blog \
6 Fun stuff My Books My Open Source Projects Hire Me Free Mentoring \
7 Privacy Policy Mark Watson: AI Practitioner and Polyglot Programmer I am the author \
8 of 20+ books on Artificial Intelligence, Common Lisp, Deep Learning, Haskell, Clojur\
9 e, Java, Ruby, Hy language, and the Semantic Web. I have 55 US Patents. My customer \
10 list includes: Google, Capital One, Olive AI, CompassLabs, Disney, SAIC, Americast, \
11 PacBell, CastTV, Lutris Technology, Arctan Group, Sitescout.com, Embed.ly, and Webmi\
12 nd Corporation.
13
14 ++ h1_headers: ["Mark Watson: AI Practitioner and Polyglot Programmer", "The books t\
15 hat I have written", "Fun stuff", "Open Source", "Hire Me", "Free Mentoring", "Priva\
16 cy Policy"]
17
18 ++ h2_headers: ["I am the author of 20+ books on Artificial Intelligence, Common Lis\
19 p, Deep Learning, Haskell, Clojure, Java, Ruby, Hy language, and the Semantic Web. I\
20 have 55 US Patents.", "Other published books:"]
21
22 ++ anchors: [["Read my Blog", "https://mark-watson.blogspot.com"], ["Fun stuff", "ht\
23 ps://markwatson.com#fun"], ["My Books", "https://markwatson.com#books"], ["My Open \
24 Source Projects", "https://markwatson.com#opensource"], ["Hire Me", "https://markwat\
25 son.com#consulting"], ["Free Mentoring", "https://markwatson.com#mentoring"], ["Priv\
26 acy Policy", "https://markwatson.com/privacy.html"], ["leanpub", "https://leanpub.co\
27 m/u/markwatson"], ["GitHub", "https://github.com/mark-watson"], ["LinkedIn", "https:\
28 //www.linkedin.com/in/marklwatson/"], ["Twitter", "https://twitter.com/mark_l_watson\
29 "], ["leanpub", "https://leanpub.com/lovinglisp"], ["leanpub", "https://leanpub.com/\\
30 haskell-cookbook/"], ["leanpub", "https://leanpub.com/javaai"],
31]
32 Test Suite 'All tests' passed at 2021-08-06 17:37:11.062.
33 Executed 2 tests, with 0 failures (0 unexpected) in 0.471 (0.472) seconds
```

## Running in the Swift REPL

```
1 $ swift run --repl
2 [1/1] Build complete!
3 Launching Swift REPL with arguments: -I/Users/markw_1/GIT_swift_book/WebScraping_swif\
4 ft/.build/arm64-apple-macosx/debug -L/Users/markw_1/GIT_swift_book/WebScraping_swift\
5 /.build/arm64-apple-macosx/debug -lWebScraping_swift__REPL
6 Welcome to Apple Swift version 5.5 (swiftlang-1300.0.29.102 clang-1300.0.28.1).
7 Type :help for assistance.
8 1> import WebScraping_swift
9 2> webPageText(uri: "https://markwatson.com")
10 $R0: String = "Mark Watson: AI Practitioner and Polyglot Programmer | Mark Watson \
11 Read my Blog Fun stuff My Books My Open Source Projects Privacy Policy \
12 Mark Watson: AI Practitioner and Polyglot Programmer I am the author of 20+ books on\
13 Artificial Intelligence, Common Lisp, Deep Learning, Haskell, Clojure, Java, Ruby, \
14 Hy language, and the Semantic Web. I have 55 US Patents. My customer list includes: \
15 Google, Capital One, Babylint, Olive AI, CompassLabs, Disney, SAIC, Americast, PacBe\
16 ll, CastTV, Lutris Technology, Arctan Group, Sitescout.com, Embed.ly, and Webmind Co\
17 rporation"...
18 3>
```

This chapter finishes a quick introduction to using Swift and Swift packages for command line utilities. The remainder of this book comprises machine learning, natural language processing, and semantic web/linked data examples.

# Part 2: Large Language Models

In this part we cover:

- Commercial OpenAI LLM APIs
- Commercial Anthropic LLM APIs
- Accessing open weight models using the commercial Groq service
- Accesing xAIs Grok model via an API
- Accessing local LLMs using Ollama
- Using Local LLMs with Apple's MLX Framework

# Using the OpenAI LLM APIs

I have been working as an artificial intelligence practitioner since 1982 and the capability of Large Language Models (LLMs) is unlike anything I have seen before. I managed a deep learning team at Capital One in 2017-2019 and we used precursors of TransFormer models like OpenAI's ChatGPT, and Anthropic's Claude.

You will need to apply to OpenAI for an access key at:

<https://platform.openai.com/signup>

The GitHub repository for this example is:

[https://github.com/mark-watson/OpenAI\\_swift](https://github.com/mark-watson/OpenAI_swift)

I recommend reading the online documentation for the [online documentation for the APIs<sup>29</sup>](https://openai.com/docs/) to see all the capabilities of the beta OpenAI APIs. Let's start by jumping into the example code that is a GitHub repository [https://github.com/mark-watson/OpenAI\\_swift<sup>30</sup>](https://github.com/mark-watson/OpenAI_swift<sup>30</sup>) that you can use in your projects.

The library that I wrote for this chapter supports four functions: for completing text, summarizing text, answering general questions, and getting embeddings for text. The get-4o-mini that we will use here is very inexpensive and capable.

You need to request an API key (I had to wait a few weeks to receive my key) and set the value of the environment variable **OPENAI\_KEY** to your key. You can add a statement like:

```
export OPENAI_KEY=sa-hdedds7&dhdhsdffd...
```

to your **.profile** or other shell resource file that contains your key value (the above key value is made-up and invalid).

The file **Sources/OpenAI\_swift/OpenAI\_swift.swift** contains the source code (code description follows the listing):

---

<sup>29</sup><https://openai.com/docs/>

<sup>30</sup>[https://github.com/mark-watson/OpenAI\\_swift](https://github.com/mark-watson/OpenAI_swift)

```
import Foundation

struct OpenAI {
 private static let key = ProcessInfo.processInfo.environment["OPENAI_KEY"]!
 private static let baseURL = "https://api.openai.com/v1"

 private struct ChatRequest: Encodable {
 let model: String
 let messages: [[String: String]]
 let max_tokens: Int
 let temperature: Double
 }

 private struct EmbeddingRequest: Encodable {
 let model: String
 let input: String
 }

 private static func makeRequest<T: Encodable>(endpoint: String, body: T)
 -> String {
 var responseString = ""
 let url = URL(string: baseURL + endpoint)!
 var request = URLRequest(url: url)
 request.httpMethod = "POST"
 request.setValue("application/json", forHTTPHeaderField: "Content-Type")
 request.setValue("Bearer \(key)", forHTTPHeaderField: "Authorization")
 request.httpBody = try? JSONEncoder().encode(body)

 let semaphore = DispatchSemaphore(value: 0)
 URLSession.shared.dataTask(with: request) { data, response, error in
 if let error = error {
 print("Error: \(error)")
 }
 if let data = data {
 responseString = String(data: data, encoding: .utf8) ?? "{}"
 }
 semaphore.signal()
 }.resume()
 semaphore.wait()

 return responseString
 }
}
```

```

static func chat(messages: [[String: String]], maxTokens: Int = 25,
 temperature: Double = 0.3) -> String {
 let chatRequest = ChatRequest(
 model: "gpt-4o-mini",
 messages: messages,
 max_tokens: maxTokens,
 temperature: temperature
)

 let response = makeRequest(endpoint: "/chat/completions", body: chatRequest)
 guard let data = response.data(using: .utf8),
 let json = try? JSONSerialization.jsonObject(with: data)
 as? [String: Any],
 let choices = json["choices"] as? [[String: Any]],
 let firstChoice = choices.first,
 let message = firstChoice["message"] as? [String: Any],
 let content = message["content"] as? String else {
 return ""
 }
 return content
}

static func embeddings(text: String) -> [Float] {
 let embeddingRequest = EmbeddingRequest(
 model: "text-embedding-ada-002",
 input: text
)

 let response = makeRequest(endpoint: "/embeddings", body: embeddingRequest)
 guard let data = response.data(using: .utf8),
 let json = try? JSONSerialization.jsonObject(with: data)
 as? [String: Any],
 let dataArray = json["data"] as? [[String: Any]],
 let embedding = dataArray.first?["embedding"] as? [NSNumber] else {
 return [1.23]
 }
 return embedding.map { number in Float(truncating: number) }
}

// Usage functions:
func summarize(text: String, maxTokens: Int = 40) -> String {
 OpenAI.chat(messages: [

```

```

 ["role": "system",
 "content":
 "You are a helpful assistant that summarizes text concisely."],
 ["role": "user", "content": text]
], maxTokens: maxTokens)
}

func questionAnswering(question: String) -> String {
 OpenAI.chat(messages: [
 ["role": "system",
 "content":
 "You are a helpful assistant that answers questions directly and concisely."],
 ["role": "user", "content": question]
], maxTokens: 25)
}

func completions(promptText: String, maxTokens: Int = 25) -> String {
 OpenAI.chat(messages: [["role": "user", "content": promptText]],
 maxTokens: maxTokens)
}

```

This Swift implementation provides a streamlined interface to OpenAI's API services, focusing primarily on chat completions and text embeddings functionality. The code is structured around a central OpenAI struct that encapsulates all API interactions and provides a clean, type-safe interface for making requests.

## Core Architecture

The implementation follows a modular design pattern, separating concerns between network communication, request/response handling, and utility functions. It utilizes Swift's strong type system through dedicated request models and leverages environment variables for secure API key management.

## Key Features

### Authentication and Configuration

The client automatically retrieves the OpenAI API key from environment variables, providing a secure way to handle authentication credentials. The base URL is configured as a constant, making it easy to modify for different environments or API versions.

## Chat Completions

The chat completion functionality supports the GPT-4 model family, allowing for structured conversations through an array of messages. Each message contains a role (system, user, or assistant) and content. The implementation provides fine-grained control over:

- Maximum token output
- Temperature settings for response randomness
- Message context management
- Text embeddings

The embeddings feature implements OpenAI's text-embedding-ada-002 model, converting text inputs into high-dimensional vector representations. These embeddings can be used for:

- Semantic search
- Text similarity comparisons
- Document classification
- Other natural language processing tasks

## Utility Functions

The implementation includes pre-built utility functions for common use cases:

- Text summarization with customizable length
- Question-answering with concise responses
- General text completions

## Technical Implementation Details

### Network Communication

The networking layer uses URLSession with a synchronous approach via DispatchSemaphore. While this ensures straightforward usage, it's worth noting that this approach should be carefully considered for production environments where asynchronous communication might be more appropriate.

### Error Handling

The implementation includes basic error handling through Swift's optional binding and guard statements, providing graceful fallbacks for common failure scenarios. The embedding function, for instance, returns a default value rather than throwing an error when processing fails.

## Data Parsing

JSON parsing is handled through a combination of JSONEncoder for requests and JSONSerialization for responses, with careful optional chaining to safely handle malformed or unexpected responses.

## Running Tests

The file SWIFT\_BOOK/OpenAI\_swift/Tests/OpenAI\_swiftTests/OpenAI\_swiftTests.swift contains test code:

```

1 import XCTest
2 @testable import OpenAI_swift
3
4 final class OpenAI_swiftTests: XCTestCase {
5 func testExample() {
6 print("Starting tests...")
7 let embeds = OpenAI.embeddings(text: "Congress passed tax laws.")
8 print(embeds[..

```

Output from this test code is:

```
$ swift test
Building for debugging...
[4/4] Compiling OpenAI_swift OpenAI_swift.swift
Build complete! (0.59s)
Test Suite 'All tests' started at 2024-11-17 16:42:12.354.
Test Suite 'OpenAI_swiftPackageTests.xctest' started at 2024-11-17 16:42:12.355.
Test Suite 'OpenAI_swiftTests' started at 2024-11-17 16:42:12.355.
Test Case '-[OpenAI_swiftTests.OpenAI_swiftTests testExample]' started.
Test Case '-[OpenAI_swiftTests.OpenAI_swiftTests testExample]' passed (7.429 seconds\).
Test Suite 'OpenAI_swiftTests' passed at 2024-11-17 16:42:19.784.
 Executed 1 test, with 0 failures (0 unexpected) in 7.429 (7.429) seconds
Test Suite 'OpenAI_swiftPackageTests.xctest' passed at 2024-11-17 16:42:19.785.
 Executed 1 test, with 0 failures (0 unexpected) in 7.429 (7.430) seconds
Test Suite 'All tests' passed at 2024-11-17 16:42:19.785.
 Executed 1 test, with 0 failures (0 unexpected) in 7.429 (7.431) seconds
Starting tests...
[-0.0046315026, -0.0077434415, 0.0005571732, -0.024781546, -0.0031119392, -0.0185893\]
25, -0.003362034, -0.020906659, 0.0074585234, -0.019463073]
** ret from OpenAI API call: the shimmering surface of the water, where the sunlight\
danced in tiny sparkles. The gentle flow of the river whispered secrets as
** answer from OpenAI API call: Leonardo da Vinci was born in Vinci, Italy, on April\
15, 1452.
** generated summary: Jupiter is the fifth planet from the Sun and the largest in t\
he Solar System, being a gas giant with a mass one-thousandth that of the Sun and tw\
o-and-a-half times that of
* Test run started.
* Testing Library Version: 102 (arm64e-apple-macos13.0)
* Test run with 0 tests passed after 0.001 seconds.
```

# Using APIs for Anthropic Claude LLMs

Here, I decided to not write a new client library for the Anthropic APIs since there are several existing high quality libraries for accessing the Anthropic Claude APIs.

This is not a strong recommendation of one Anthropic client library over another, but I very much enjoy using the following project because of the simplicity of its API:

<https://github.com/fumito-ito/AnthropicSwiftSDK>

My examples using this library to access the Anthropic Claude APIs can be found here:

[https://github.com/mark-watson/Anthropic\\_swift\\_examples](https://github.com/mark-watson/Anthropic_swift_examples)

You need to set the following environment variable for your person Anthropic API key: Anthropic API key:

ANTHROPIC\_API\_KEY

that you can get by creating an account:

<https://console.anthropic.com>

Note that there is no library implemented in this chapter.

## Running the examples

All of the examples are packaged as Swift tests so git clone my examples repository [https://github.com/mark-watson/Anthropic\\_swift\\_examples](https://github.com/mark-watson/Anthropic_swift_examples)<sup>31</sup> and run:

`swift test`

The test Swift source file defines a test class (just the first few lines shown here):

---

<sup>31</sup>[https://github.com/mark-watson/Anthropic\\_swift\\_examples](https://github.com/mark-watson/Anthropic_swift_examples)

```

final class Anthropic_swift_examplesTests: XCTestCase {
 let text1 = "If Mary is 42, Bill is 27, and Sam is 51, what are their pairwise age\
differences. Please be concise."
 func testExample() async throws {
 let key =
 ProcessInfo.processInfo.environment["ANTHROPIC_API_KEY"]!
 let anthropic = Anthropic(apiKey: key)
 let message = Message(role: .user,
 content: [.text(text1)])
 let response =
 try await
 anthropic.messages.createMessage([message],
 maxTokens: 400)
 }
}

```

If you print the value of **response** you see:

```

MessagesResponse(id: "msg_02RvimFjHFFmaV4n994J9wck", type: AnthropicSwiftSDK.Message\
sResponseType.message, role: AnthropicSwiftSDK.Role.assistant, content: [AnthropicSwi\
ftSDK.Content.text("The pairwise age differences are:\n\nMary and Bill: 15 years\nM\
ary and Sam: 9 years\nBill and Sam: 24 years", cacheControl: nil)], model: Optional(\\
AnthropicSwiftSDK.Model.claude_3_Opus), stopReason: Optional(AnthropicSwiftSDK.StopR\
eason.endTurn), stopSequence: nil, usage: AnthropicSwiftSDK.TokenUsage(inputTokens: \
Optional(38), outputTokens: Optional(38)))

```

If you print the value of **response.content** you see:

```
[AnthropicSwiftSDK.Content.text("The pairwise age differences are:\n\nMary and Bill:\
15 years\nMary and Sam: 9 years\nBill and Sam: 24 years", cacheControl: nil)]
```

For normal use you want just the string contents of the model's response to your prompt, so use:

```

for content in response.content {
 if case let .text(text, _) = content {
 print("Assistant's response: \(text)")
 }
}

```

That outputs:

Assistant's response: The pairwise age differences are:

Mary and Bill: 15 years

Mary and Sam: 9 years

Bill and Sam: 24 years

In the general case of the Claude model returning images, tools used, and tool results, use code like this:

```
for content in response.content {
 switch content {
 case .text(let text, _):
 print("Assistant's response: \(text)")
 case .image(let imageContent, _):
 print("imageContent: \(imageContent)")
 break
 case .document(let documentContent, _):
 print("documentContent: \(documentContent)")
 break
 case .toolResult(let toolResult):
 print("toolResult: \(toolResult)")
 break
 case .toolUse(let toolUse):
 print("toolUse: \(toolUse)")
 break
 }
}
```

# Using Groq APIs to Open Weight LLM Models

Groq develops custom silicon for fast LLM inference.

Groq's API service supports a variety of openly available models, including:

- Llama 3.1 Series: Models like llama-3.1-70b-versatile, llama-3.1-8b-instant, and others, offering up to 128K context windows.
- Llama 3.2 Vision Series: Multimodal models such as llama-3.2-90b-vision-preview and llama-3.2-11b-vision-preview, capable of processing both text and image inputs.
- Llama 3 Groq Tool Use Models: Specialized for function calling, including llama3-Groq-70b-8192-tool-use-preview and llama3-Groq-8b-8192-tool-use-preview.
- Mixtral 8x7b: A model with a 32,768-token context window, suitable for extensive context applications.
- Gemma Series: Models like gemma2-9b-it and gemma-7b-it, each with an 8,192-token context window.
- Whisper Series: Models such as whisper-large-v3 and whisper-large-v3-turbo, designed for audio transcription and translation tasks.

To obtain an API key, visit Groq's API keys management page:

<https://console.groq.com/keys>

The code for this chapter can be found here:

[https://github.com/mark-watson/Groq\\_swift](https://github.com/mark-watson/Groq_swift)

## Implementation of a Client Library for the Groq APIs

Groq supports the OpenAI APIs so the following client library for Groq is similar to what I wrote previously for OpenAI:

```

import Foundation

struct Groq {
 private static let key = ProcessInfo.processInfo.environment["GROQ_API_KEY"]!
 private static let baseURL = "https://api.groq.com/openai/v1/"

 private static let MODEL = "llama3-8b-8192"

 private struct ChatRequest: Encodable {
 let model: String
 let messages: [[String: String]]
 let max_tokens: Int
 let temperature: Double
 }

 private static func makeRequest<T: Encodable>(endpoint: String, body: T)
 -> String {
 var responseString = ""
 let url = URL(string: baseURL + endpoint)!
 var request = URLRequest(url: url)
 request.httpMethod = "POST"
 request.setValue("application/json", forHTTPHeaderField: "Content-Type")
 request.setValue("Bearer \(key)", forHTTPHeaderField: "Authorization")
 request.httpBody = try? JSONEncoder().encode(body)

 let semaphore = DispatchSemaphore(value: 0)
 URLSession.shared.dataTask(with: request) { data, response, error in
 if let error = error {
 print("Error: \(error)")
 }
 if let data = data {
 responseString = String(data: data, encoding: .utf8) ?? "{}"
 }
 semaphore.signal()
 }.resume()
 semaphore.wait()

 return responseString
 }

 static func chat(messages: [[String: String]],
 maxTokens: Int = 25,
 temperature: Double = 0.3)

```

```

 -> String {
 let chatRequest = ChatRequest(
 model: MODEL,
 messages: messages,
 max_tokens: maxTokens,
 temperature: temperature
)

 let response = makeRequest(endpoint: "/chat/completions", body: chatRequest)
 guard let data = response.data(using: .utf8),
 let json = try? JSONSerialization.jsonObject(with: data)
 as? [String: Any],
 let choices = json["choices"] as? [[String: Any]],
 let firstChoice = choices.first,
 let message = firstChoice["message"]
 as? [String: Any],
 let content = message["content"] as? String else {
 return ""
 }
 return content
 }

// Usage functions:
func summarize(text: String, maxTokens: Int = 40) -> String {
 Groq.chat(messages: [
 ["role": "system",
 "content": "You are a helpful assistant that summarizes text concisely"],
 ["role": "user", "content": text]
], maxTokens: maxTokens)
}

func questionAnswering(question: String) -> String {
 Groq.chat(messages: [
 ["role": "system",
 "content": "You are a helpful assistant that answers questions directly and concisely."],
 ["role": "user", "content": question]
], maxTokens: 25)
}

```

```
func completions(promptText: String, maxTokens: Int = 25) -> String {
 Groq.chat(messages: [
 ["role": "system", "content": "You complete text"],
 ["role": "user", "content": promptText]],
 maxTokens: maxTokens)
}
```

## Explanation of the Swift Groq API Code

### 1. Setting Up the API

- The `Groq` struct is designed to interact with the Groq API, mimicking OpenAI's API.
- **API Key:** Retrieved from the environment variable `GROQ_API_KEY`.
- **Base URL:** The API's base endpoint is `https://api.groq.com/openai/v1/`.
- **Model:** The model being used is predefined as `11ama3-8b-8192`.

### 2. Structure of a Chat Request

- A private struct, `ChatRequest`, defines the JSON payload for requests:
    - `model`: The model name.
    - `messages`: A history of the conversation.
    - `max_tokens`: Limits the number of tokens in the response.
    - `temperature`: Controls randomness in the responses.
- 

### 3. Making an HTTP POST Request

- The `makeRequest` function handles API communication:
    - Constructs the full URL by appending the endpoint to the base URL.
    - Sets up a POST request with:
      - \* JSON content type.
      - \* Authorization header using the API key.
    - Encodes the request body into JSON using `JSONEncoder`.
    - Sends the request asynchronously but waits for the response using a semaphore.
    - Parses the response data into a string.
-

## 4. Chat Functionality

- The chat function simplifies sending messages to the API:
    - Constructs a `ChatRequest` object with the given parameters.
    - Sends the request to the `/chat/completions` endpoint.
    - Processes the JSON response to extract the model’s reply from the `choices` array.
- 

## 5. Usage Functions

**summarize** - Summarizes a text. - Sends a conversation history where the system is described as “a helpful assistant that summarizes text concisely.”

```
1 summarize(text: String, maxTokens: Int = 40)
```

**questionAnswering** - Answers a user-provided question directly. - Sends a conversation history where the system is described as “a helpful assistant that answers questions directly and concisely.”

**completions** - Generates continuations for a given user prompt.

## Running the Tests

Here is the test/example code for this library:

```
import XCTest
@testable import Groq_swift

final class Groq_swiftTests: XCTestCase {
 func testExample() {
 print("Starting tests...")
 let prompt = "He walked to the river and looked at"
 let ret = completions(promptText: prompt, maxTokens: 200)
 print("** ret from Groq API call:", ret)
 let question = "Where was Leonardo da Vinci born?"
 let answer = questionAnswering(question: question)
 print("** answer from Groq API call:", answer)
 let text = "Jupiter is the fifth planet from the Sun and the largest in the Solar System. It is a gas giant with a mass one-thousandth that of the Sun, but two-and-a-half times that of all the other planets in the Solar System combined. Jupiter is one of the brightest objects visible to the naked eye in the night sky, and has been known to ancient civilizations since before recorded history. It is named after the\
```

Roman god Jupiter.[19] When viewed from Earth, Jupiter can be bright enough for its reflected light to cast visible shadows,[20] and is on average the third-brightest natural object in the night sky after the Moon and Venus."

```
let summary = summarize(text: text)
print("** generated summary: ", summary)
}
}
```

Here is sample output from this example use of the library:

Starting tests...

```
** ret from Groq API call: the calm water, feeling the warm sun on his face and the \
gentle breeze rustling his hair. The sound of the water lapping against the shore wa\
s soothing, and he closed his eyes, taking a deep breath to clear his mind.
```

```
** answer from Groq API call: Leonardo da Vinci was born on April 15, 1452, in Vinci\
, Italy.
```

```
** generated summary: Here's a concise summary:
```

```
Jupiter is the largest planet in our Solar System, a gas giant with a mass 2.5 times\
that of all other planets combined. It's the fifth planet
```

# Using the xAI Grok LLM

xAI's Grok is a large language model (LLM) developed by Elon Musk's AI startup, xAI, to compete with leading AI systems like OpenAI's GPT-4. Launched in 2023, Grok is designed to handle a variety of tasks, including answering questions, assisting with writing, and solving coding problems. It is integrated with the social media platform X (formerly Twitter), providing users with real-time information and a conversational AI experience.

Grok has undergone several iterations, with Grok-2 being released in August 2024. This version introduced image generation capabilities, enhancing its versatility. xAI has also made Grok-1 open-source, allowing developers to access its weights and architecture for further research and application development.

To support Grok's development, xAI has invested in substantial computational resources, including the Colossus supercomputer, which utilizes 100,000 Nvidia H100 GPUs, positioning it as one of the most powerful AI training systems globally.

## Implementation of a Grok API Client Library

The code for my xAI Grok client code can be found here:

[https://github.com/mark-watson/X\\_GROK\\_swift](https://github.com/mark-watson/X_GROK_swift)

The Grok is similar to the OpenAI APIs so I copied the code we saw earlier that I wrote for OpenAI and made the simple modifications required to access Grok APIs (code discussion appears after this listing):

```
import Foundation

// xAI Grok LLM client library

struct X_GROK {
 private static let key = ProcessInfo.processInfo.environment["X_GROK_API_KEY"]!
 private static let baseURL = "https://api.x.ai/v1"

 private static let MODEL = "grok-beta"

 private struct ChatRequest: Encodable {
 let model: String
```

```
let messages: [[String: String]]
let max_tokens: Int
let temperature: Double
}

private static func makeRequest<T: Encodable>(endpoint: String, body: T)
 -> String {
 var responseString = ""
 let url = URL(string: baseURL + endpoint)!
 var request = URLRequest(url: url)
 request.httpMethod = "POST"
 request.setValue("application/json", forHTTPHeaderField: "Content-Type")
 request.setValue("Bearer \(key)", forHTTPHeaderField: "Authorization")
 request.httpBody = try? JSONEncoder().encode(body)

 let semaphore = DispatchSemaphore(value: 0)
 URLSession.shared.dataTask(with: request) { data, response, error in
 if let error = error {
 print("Error: \(error)")
 }
 if let data = data {
 responseString = String(data: data, encoding: .utf8) ?? "{}"
 }
 semaphore.signal()
 }.resume()
 semaphore.wait()

 return responseString
}

static func chat(messages: [[String: String]], maxTokens: Int = 25,
 temperature: Double = 0.3) -> String {
 let chatRequest = ChatRequest(
 model: MODEL,
 messages: messages,
 max_tokens: maxTokens,
 temperature: temperature
)

 let response = makeRequest(endpoint: "/chat/completions",
 body: chatRequest)
 guard let data = response.data(using: .utf8),
 let json = try? JSONSerialization.jsonObject(with: data) as? [String: \
```

```

Any] ,
 let choices = json["choices"] as? [[String: Any]],
 let firstChoice = choices.first,
 let message = firstChoice["message"]
 as? [String: Any],
 let content = message["content"] as? String else {
 return ""
 }
 return content
}
}

// Usage functions:
func summarize(text: String, maxTokens: Int = 40) -> String {
 X_GROK.chat(messages: [
 ["role": "system",
 "content":
 "You are a helpful assistant that summarizes text concisely."],
 ["role": "user", "content": text]
], maxTokens: maxTokens)
}

func questionAnswering(question: String) -> String {
 X_GROK.chat(messages: [
 ["role": "system",
 "content":
 "You are a helpful assistant who answers questions directly and concisely."],
 ["role": "user", "content": question]
], maxTokens: 25)
}

func completions(promptText: String, maxTokens: Int = 25) -> String {
 X_GROK.chat(messages: [[{"role": "user", "content": promptText}],
 maxTokens: maxTokens)
}

```

This Swift code defines a client library, `X_GROK`, to interact with the xAI Grok Large Language Model (LLM) API. It leverages Swift's Foundation framework to handle HTTP requests and JSON encoding/decoding. The library retrieves the API key from the environment variable `X_GROK_API_KEY` and sets the base URL for the API. It specifies a default model, `grok-beta`, for generating responses.

The core functionality is encapsulated in the `makeRequest` function, which constructs and sends HTTP POST requests to the API. It accepts an endpoint and a request body conforming to the

Encodable protocol. The function sets the necessary HTTP headers, including Content-Type and Authorization, and encodes the request body into JSON. To handle the asynchronous nature of network calls synchronously, it employs a semaphore, ensuring the function waits for the response before proceeding. The response is then returned as a string.

The chat function utilizes makeRequest to send chat messages to the API. It constructs a ChatRequest struct with parameters like the model, messages, maximum tokens, and temperature. After receiving the response, it parses the JSON to extract the generated content. Additionally, the code provides utility functions **summarize**, **questionAnswering**, and **completions** which use the **chat** function to perform specific tasks such as text summarization, question answering, and text completion, respectively.

Here is the test/example code for this library:

```
import XCTest
@testable import X_GROK_swift

final class X_GROK_swiftTests: XCTestCase {
 func testExample() {
 print("Starting tests...")
 let prompt = "He walked to the river and looked at"
 let ret = completions(promptText: prompt, maxTokens: 200)
 print("** ret from X_GROK API call:", ret)
 let question = "Where was Leonardo da Vinci born?"
 let answer = questionAnswering(question: question)
 print("** answer from X_GROK API call:", answer)
 let text = "Jupiter is the fifth planet from the Sun and the largest in the Solar System. It is a gas giant with a mass one-thousandth that of the Sun, but two-and-a-half times that of all the other planets in the Solar System combined. Jupiter is one of the brightest objects visible to the naked eye in the night sky, and has been known to ancient civilizations since before recorded history. It is named after the Roman god Jupiter.[19] When viewed from Earth, Jupiter can be bright enough for its reflected light to cast visible shadows,[20] and is on average the third-brightest natural object in the night sky after the Moon and Venus."
 let summary = summarize(text: text)
 print("** generated summary: ", summary)
 }
}
```

Here is the example code output:

Starting tests...

\*\* ret from X\_GROK API **call**: He walked to the river **and** looked at the water flowing \ gently by. The serene scene provided a moment of peace, as he watched the ripples da\ nce in the sunlight. The sound of the water was soothing, almost like a lullaby, cal\ ming his thoughts **and** grounding him in the present. He felt a deep connection to nat\ ure, the river's timeless flow reminding him of life's continuous journey.

\*\* answer from X\_GROK API **call**: Leonardo da Vinci was born in the town of Vinci, in \ the region of Tuscany, Italy.

\*\* generated **summary**: Jupiter, the fifth planet from the Sun, is the largest in our\ Solar System, classified as a gas giant with a mass significantly greater than all \ other planets combined. It's one of the brightest objects

# Using Ollama to Run Local LLMs

Ollama is a program and framework written in Go that allows you to download, run models on the command line, and call using a REST style interface. You need to download the Ollama executable for your operation system at <https://ollama.com><sup>32</sup>.

Similarly to our use of a third party for accessing the Anthropic Claude models, here we will not write a wrapper library. The example code for this chapter is in the test code for the Swift project in the GitHub repository [https://github.com/mark-watson/Ollama\\_swift\\_examples](https://github.com/mark-watson/Ollama_swift_examples)<sup>33</sup>.

We use the library in the GitHub repository <https://github.com/mattt/ollama-swift><sup>34</sup>.

## Running the Ollama Service

Assuming you have Ollama installed, download the following model that required two gigabytes of disk space:

```
ollama pull llama3.2:latest
```

When the model is downloaded it is also cached for future use on your laptop.

Here is the test/example code we will run:

```
import XCTest
import Ollama

final class Ollama_swift_examplesTests: XCTestCase {
 let text1 = "If Mary is 42, Bill is 27, and Sam is 51, what are their pairwise age differences."
 let client = Ollama.Client.default // http://localhost:11434 endpoint
 func testExample() async throws {
 let response = try await client.chat(
 model: "llama3.2:latest",
 messages: [
 .system("You are a helpful assistant who completes text and also answers questions. You are always concise."),
 .user(text1),
```

---

<sup>32</sup><https://ollama.com>

<sup>33</sup>[https://github.com/mark-watson/Ollama\\_swift\\_examples](https://github.com/mark-watson/Ollama_swift_examples)

<sup>34</sup><https://github.com/mattt/ollama-swift>

```

 .user("what if Sam is 52?")
])
print(response.message.content)
}
}
```

The output looks like:

Pairwise age differences:

- Mary - Bill:  $|42 - 27| = 15$
- Mary - Sam:  $|42 - 51| = 9$
- Bill - Sam:  $|27 - 51| = 24$

If Sam is 52:

- Mary - Bill:  $|42 - 27| = 15$
- Mary - Sam:  $|42 - 52| = 10$
- Bill - Sam:  $|27 - 52| = 25$

The **ollama\_swift** library also supports text generation. You can also do single shot text generation using the code in the previous example, but only using one **user** call, for example:

```

final class Ollama_swift_examplesTests: XCTestCase {
 let text1 = "What is the capital of Germany?"
 let client = Ollama.Client.default
 func testExample() async throws {
 let response = try await client.chat(
 model: "llama3.2:latest",
 messages: [
 .system("You are a helpful assistant who completes text and also answers\
questions. You are always concise."),
 .user(text1),
])
 print(response.message.content)
 }
}
```

The output looks like:

The capital of Germany is Berlin.

## Ollama Wrap Up

This is a short chapter but an important one. I do over half my work with LLMs running locally on my laptop using Ollama, with the rest of my work using OpenAI, Anthropic, and Groq commercial APIs.

# Using Apple's MLX Framework to Run Local LLMs

Apple's MLX framework is an efficient way to use LLMs embedded in applications written in Swift using the SwiftUI user interface library for macOS, iOS, and iPadOS.

It is difficult to create simple command line Swift apps using MLX but there are several complete MLX, Swift, and SwiftUI demo applications that you can use to start your own projects. Here we will use the LLMEval application from the GitHub repository [https://github.com/ml-explore/mlx-swift-examples<sup>35</sup>](https://github.com/ml-explore/mlx-swift-examples).

## MLX Framework History

Apple's MLX framework, introduced in December 2023, is a key part of Apple's strategy to support AI on its hardware platforms by leveraging the unique capabilities of Apple Silicon, including the M1, M2, M3, and M4 series. Designed as an open-source, NumPy-like array framework, MLX optimizes machine learning workloads, particularly large language models (LLMs), by utilizing Apple Silicon's unified architecture that integrates CPU, GPU, Neural Engine, and shared memory. This architecture eliminates data transfer bottlenecks, enabling faster and more efficient ML tasks, such as training and deploying LLMs directly on devices like MacBooks and iPhones. MLX aligns with Apple's privacy-focused approach by supporting on-device processing, enhancing performance for applications like natural language processing, speech recognition, and content generation while offering a seamless transition for Python or Swift ML engineers familiar with frameworks like NumPy and PyTorch. MLX stands out by leveraging Apple's unified memory architecture, allowing shared memory access between CPU and GPU, which eliminates data transfer overhead and accelerates machine learning tasks, especially with large datasets.

## MLX Resources on GitHub

In this chapter we will look at an example application that is part of the Swift MLX Examples project. After working through this example, the following resources on GitHub are worth looking at:

- <https://github.com/ml-explore/mlx-swift>: The Swift API for MLX, enabling integration with Swift-based projects.

---

<sup>35</sup><https://github.com/ml-explore/mlx-swift-examples>

- <https://github.com/ml-explore/mlx-swift-examples>: Examples showcasing the use of MLX with Swift.

You can find the documentation here:

<https://swiftpackageindex.com/ml-explore/mlx-swift/0.18.0/documentation/mlx<sup>36</sup>>.

These repositories provide a comprehensive set of tools and examples to effectively utilize MLX for machine learning tasks on Apple silicon. There are many other repositories for MLX and Python and if you need to perform tasks like fine tuning a MLX model, that task should probably be done using Python.

## Example Application for MLX Swift Examples Repository

You will want to download the complete MLX Swift examples repository:

```
git clone https://github.com/ml-explore/mlx-swift-examples.git
```

Open the top level XCode project by:

```
cd mlx-swift-examples
open mlx-swift-examples.xcodeproj
```

Here is the file browser view of this project:

---

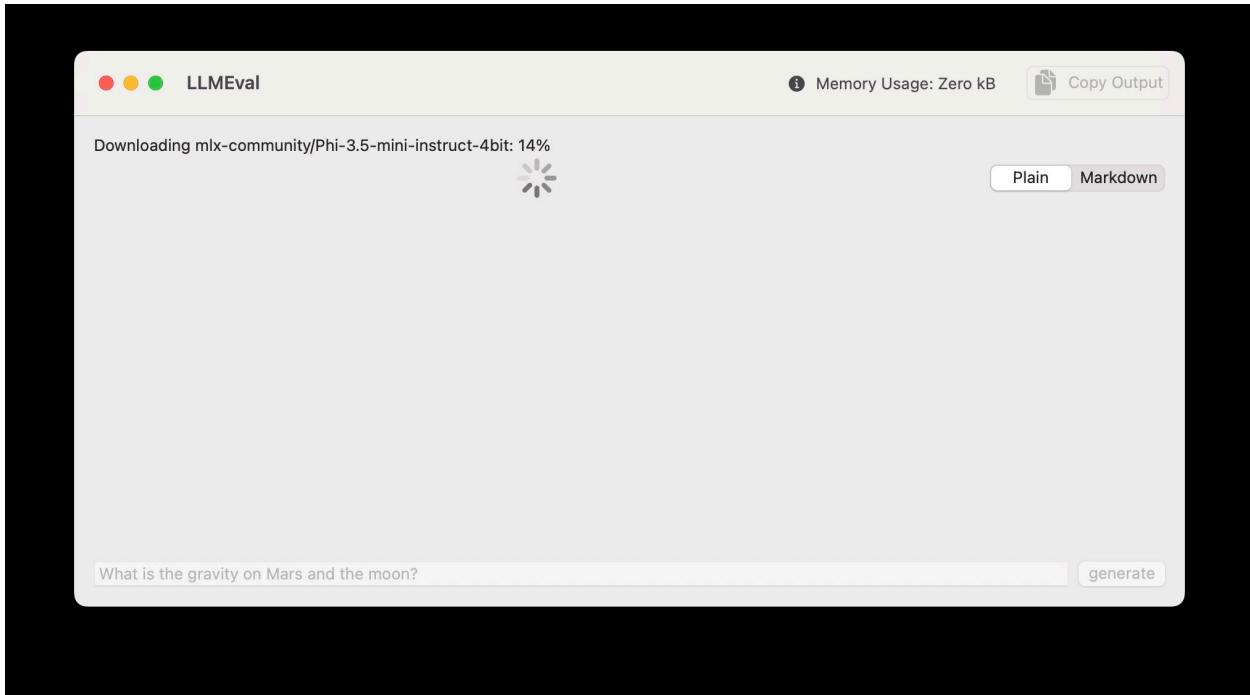
<sup>36</sup><https://swiftpackageindex.com/ml-explore/mlx-swift/0.18.0/documentation/mlx>

The screenshot shows the XCode interface with the project **mlx-swift-examples** open. The left sidebar displays the project structure:

- mlx-swift-examples**
  - README
  - Package
  - Configuration
  - Data
  - Libraries
  - Applications**
    - LoRATrainingExample
    - StableDiffusionExample
    - LLMEval**
      - ViewModels
      - Assets
      - ContentView
      - LLMEval**
      - LLMEvalApp**
    - Preview Content
    - README
    - MNISTTrainer
    - Tools
    - Products
    - Frameworks
  - Package Dependencies
    - Gzip 6.0.1
    - Ninja 1.0.5**
      - README
      - Package
      - Sources
      - Tests
      - LICENSE
    - mlx-swift 0.18.1
    - NetworkImage 6.0.0

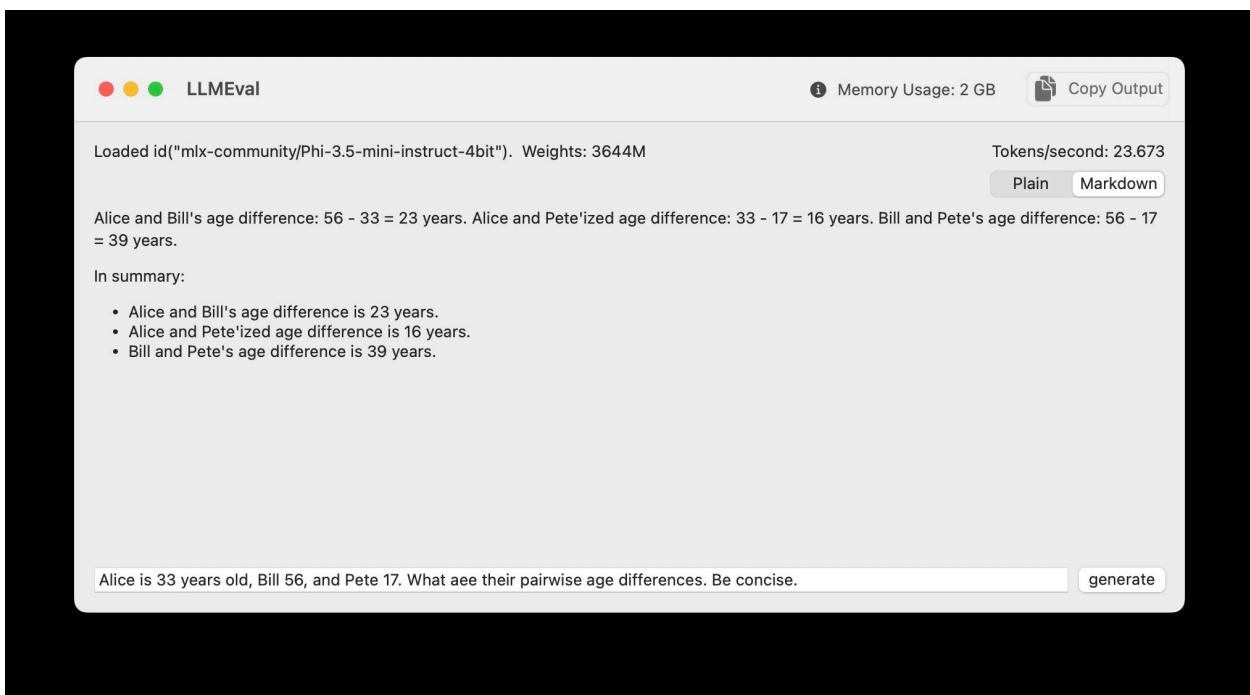
XCode View of projects

Running the LLMEval project:



LLMEval app downloading model file

Initially the model is downloaded and cached on your laptop for future use. Here is the app used to solve a simple word problem:



LLMEval app answering user's question

# Analysis of Swift and SwiftUI Code in the LLMEval Application

This example is part of the Swift MLX Examples project that currently has twenty contributors and a thousand stars on GitHub [https://github.com/ml-explore/mlx-swift-examples<sup>37</sup>](https://github.com/ml-explore/mlx-swift-examples).

Unfortunately the SwiftUI user interface code is mixed in with the code that uses MLX. Let's walk through the code:

Here is a walk through a Swift-based program using Apple's frameworks for Machine Learning and Language Models with the code interspersed with explanations.

## Imports

```
import LLM
import MLX
import MLXRandom
import MarkdownUI
import Metal
import SwiftUI
import Tokenizers
```

These imports bring in essential libraries:

- LLM and MLX for working with language models.
- MarkdownUI for rendering Markdown content.
- SwiftUI for creating the user interface.
- Tokenizers for tokenizing text.

## The ContentView Struct

The ContentView struct defines the main interface of the app.

### State Variables

```
struct ContentView: View {
 @State var prompt = ""
 @State var llm = LLMEvaluator()
 @Environment(DeviceStat.self) private var deviceStat
```

---

<sup>37</sup><https://github.com/ml-explore/mlx-swift-examples>

In this code snippet:

- `@State` allows the view to track changes in the prompt and `llm` instances.
- `@Environment` fetches device statistics, such as GPU memory usage.

## Display Style Enum

```
enum displayStyle: String, CaseIterable, Identifiable {
 case plain, markdown
 var id: Self { self }
}

@State private var selectedDisplayStyle = displayStyle.markdown
```

In this code snippet:

- `displayStyle` defines whether the output is plain text or Markdown.
- A segmented picker toggles between the two styles.

## UI Layout

### Input Section

```
var body: some View {
 VStack(alignment: .leading) {
 VStack {
 HStack {
 Text(llm.modelInfo).textFieldStyle(.roundedBorder)
 Spacer()
 Text(llm.stat)
 }
 HStack {
 Spacer()
 if llm.running {
 ProgressView().frame(maxHeight: 20)
 Spacer()
 }
 Picker("", selection: $selectedDisplayStyle) {
 ForEach(displayStyle.allCases, id: \.self) {
 option in
 Text(option.rawValue.capitalized)
 }
 }
 }
 }
 }
}
```

```

 .tag(option)
 }

 }.pickerStyle(.segmented)
}
}

```

This code displays model information and statistics.

## Output Section

```

ScrollView(.vertical) {
 ScrollViewReader { sp in
 Group {
 if selectedDisplayStyle == .plain {
 Text(llm.output)
 .textSelection(.enabled)
 } else {
 Markdown(llm.output)
 .textSelection(.enabled)
 }
 }
 .onChange(of: llm.output) { _, _ in
 sp.scrollTo("bottom")
 }
 }
}

HStack {
 TextField("prompt", text: $prompt)
 .onSubmit(generate)
 .disabled(llm.running)
 Button("generate", action: generate)
 .disabled(llm.running)
}
}

```

The ScrollView shows the model's output, which updates dynamically as the model generates text.

## Toolbar

```

.toolbar {
 ToolbarItem {
 Label(
 "Memory Usage: \(deviceStat.gpuUsage.activeMemory.formatted(.byteCount(s\
tyle: .memory)))",
 systemImage: "info.circle.fill"
)
 }
 ToolbarItem(placement: .primaryAction) {
 Button {
 Task {
 copyToClipboard(lm.output)
 }
 } label: {
 Label("Copy Output", systemImage: "doc.on.doc.fill")
 }
 }
 }
}

```

The toolbar includes:

- GPU memory usage information.
- A “Copy Output” button to copy the generated text.

## The LLMEvaluator Class

This class handles the logic for loading and generating text with the language model.

### Core Properties

```

@Observable
@MainActor
class LLMEvaluator {
 var running = false
 var output = ""
 var modelInfo = ""
 var stat = ""

 let modelConfiguration = ModelConfiguration.phi3_5_4bit
 /// parameters controlling the output
 let generateParameters = GenerateParameters(temperature: 0.6)
 let maxTokens = 240

 /// update the display every N tokens -- 4 looks like it updates continuously
}

```

```

/// and is low overhead. observed ~15% reduction in tokens/s when updating
/// on every token
let displayEveryNTokens = 4

enum LoadState {
 case idle
 case loaded(ModelContainer)
}

var loadState = LoadState.idle

```

This code snippet:

- Tracks the model state and output.
- Configures the model (phi3\_5\_4bit).

### Loading the Model (if required)

```

/// load and return the model -- can be called
/// multiple times, subsequent calls will
/// just return the loaded model

func load() async throws -> ModelContainer {
 switch loadState {
 case .idle:
 MLX.GPU.set(cacheLimit: 20 * 1024 * 1024)
 let modelContainer =
 try await LLM.loadModelContainer(configuration: modelConfiguration) {
 [modelConfiguration] progress in
 Task { @MainActor in
 self.modelInfo = "Downloading \((modelConfiguration.name): \(Int(progress.fractionCompleted * 100))%"
 }
 }
 self.modelInfo = "Loaded \((modelConfiguration.id). Weights: \(numParams / (1\(
 024*1024))M"
 loadState = .loaded(modelContainer)
 return modelContainer
 case .loaded(let modelContainer):
 return modelContainer
 }
}

```

This code snippet:

- Downloads and caches the model.
- Updates modelInfo during the download.

## Generating Output

```
func generate(prompt: String) async {
 guard !running else { return }
 running = true
 self.output = ""

 do {
 let modelContainer = try await load()
 let messages = [{"role": "user", "content": prompt}]
 let promptTokens = try await modelContainer.perform { _, tokenizer in
 try tokenizer.applyChatTemplate(messages: messages)
 }

 let result = await modelContainer.perform { model, tokenizer in
 LLM.generate(
 promptTokens: promptTokens,
 parameters: generateParameters,
 model: model,
 tokenizer: tokenizer,
 extraEOSTokens: modelConfiguration.extraEOSTokens
) { tokens in
 if tokens.count % displayEveryNTokens == 0 {
 let text = tokenizer.decode(tokens: tokens)
 Task { @MainActor in
 self.output = text
 }
 }
 if tokens.count >= maxTokens {
 return .stop
 } else {
 return .more
 }
 }
 }
 } catch {
 output = "Failed: \(error)"
 }
}
```

```
 running = false
}
}
```

This code snippet:

- Prepares the prompt for the model.
- Generates tokens and dynamically updates the view.

This program demonstrates how to integrate ML and UI components for interactive LLM-based applications in Swift.

This code example uses the MIT License so you can modify the example code if you need to write a combined SwiftUI GUI app that uses LLM-based text generation.

# **Part 3: Apple's CoreML and NLP Libraries**

In this part we cover:

- Short introduction to the ideas behind Deep Learning
- Introduction of CoreML
- Examples using CoreML
- Introduction of NLP
- Examples using NLP libraries

This section used to contain Apple CoreML examples to train a back-propagation model from the University of Wisconsin cancer data set. As of April 2022, these example do not work because of a problem with latest CreateML library so this material has been removed from this book.

# Deep Learning Introduction

Apple's work in smoothly integrating deep learning into their developer tools for macOS, iOS, and iPadOS applications is in my opinion nothing short of brilliant. We will finish this book with an application that uses two deep learning models that provide almost all of the functionality of the application.

Before diving into Apple's CoreML libraries in later chapters we will take a shallow dive into the principles of deep learning and take a lay-of-the-land look at the type of most commonly used models. This chapter has no example programs and is intended as background material.

Most of my professional career since 2014 has involved Deep Learning, mostly with TensorFlow using the Keras APIs. In the late 1980s I was on a DARPA neural network technology advisory panel for a year, I wrote the first prototype of the SAIC ANSim neural network library commercial product, and I wrote the neural network prediction code for a bomb detector my company designed and built for the FAA for deployment in airports. More recently I have used GAN (generative adversarial networks) models for synthesizing numeric spreadsheet data and LSTM (long short term memory) models to synthesize highly structured text data like nested JSON and for NLP (natural language processing). I have also written a product recommendation model for an online store using TensorFlow Recommenders. I have several USA and European patents using neural network and Deep Learning technology.

Here we will learn a vocabulary for discussing Deep Learning neural network models and look at possible architectures.

If you want to use Deep Learning professionally, there are two specific online resources that I recommend: Andrew Ng leads the efforts at [deeplearning.ai](https://deeplearning.ai)<sup>38</sup> and Jeremy Howard leads the efforts at [fast.ai](https://fast.ai)<sup>39</sup>.

There are many Deep Learning neural architectures in current practical use; a few types that I use are:

- Multi-layer perceptron networks with many fully connected layers. An input layer contains placeholders for input data. Each element in the input layer is connected by a two-dimensional weight matrix to each element in the first hidden layer. We can use any number of fully connected hidden layers, with the last hidden layer connected to an output layer.
- Convolutional networks for image processing and text classification. Convolutions, or filters, are small windows that can process input images (filters are two-dimensional) or sequences like text (filters are one-dimensional). Each filter uses a single set of learned weights independent of where the filter is applied in an input image or input sequence.

---

<sup>38</sup><https://www.deeplearning.ai/>

<sup>39</sup><https://www.fast.ai/>

- Autoencoders have the same number of input layer and output layer elements with one or more hidden fully connected layers. Autoencoders are trained to produce the same output as training input values using a relatively small number of hidden layer elements. Autoencoders are capable of removing noise in input data.
- LSTM (long short term memory) process elements in a sequence in order and are capable of remembering patterns that they have seen earlier in the sequence.
- GAN (generative adversarial networks) models comprise two different and competing neural models, the generator and the discriminator. GANs are often trained on input images (although in my work I have applied GANs to two-dimensional numeric spreadsheet data). The generator model takes as input a “latent input vector” (this is just a vector of specific size with random values) and generates a random output image. The weights of the generator model are trained to produce random images that are similar to how training images look. The discriminator model is trained to recognize if an arbitrary output image is original training data or an image created by the generator model. The generator and discriminator models are trained together.

The core functionality of libraries like TensorFlow are written in C++ and take advantage of special hardware like GPUs, custom ASICs, and devices like Google’s TPUs. Most people who work with Deep Learning models don’t need to even be aware of the low level optimizations used to make training and using Deep Learning models more efficient. That said, in the following section I am going to show you how simple neural networks are trained and used.

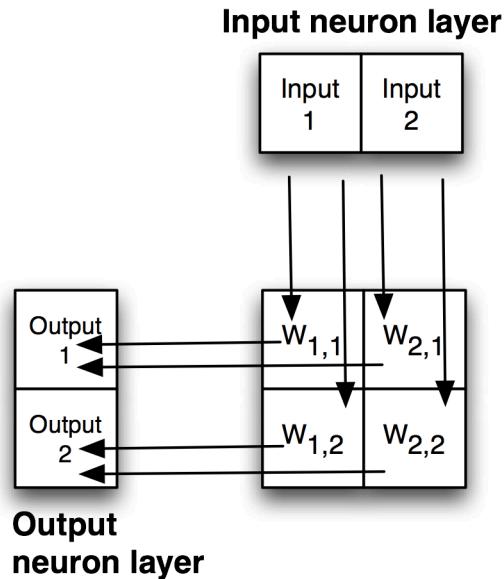
## Simple Multi-layer Perceptron Neural Networks

I use the terms Multi-layer perceptron neural networks, backpropagation neural networks and delta-rule networks interchangeably. Backpropagation refers to the model training process of calculating the output errors when training inputs are passed in the forward direction from input layer, to hidden layers, and then to the output layer. There will be an error which is the difference between the calculated outputs and the training outputs. This error can be used to adjust the weights from the last hidden layer to the output layer to reduce the error. The error is then propagated backwards through the hidden layers, updating all weights in the model. I have detailed example code in several of my older artificial intelligence books. Here I am satisfied to give you an intuition of how simple neural networks are trained.

The basic idea is that we start with a network initialized with random weights and for each training case we propagate the inputs through the network towards the output neurons, calculate the output errors, and back-up the errors from the output neurons back towards the input neurons in order to make small changes to the weights to lower the error for the current training example. We repeat this process by cycling through the training examples many times.

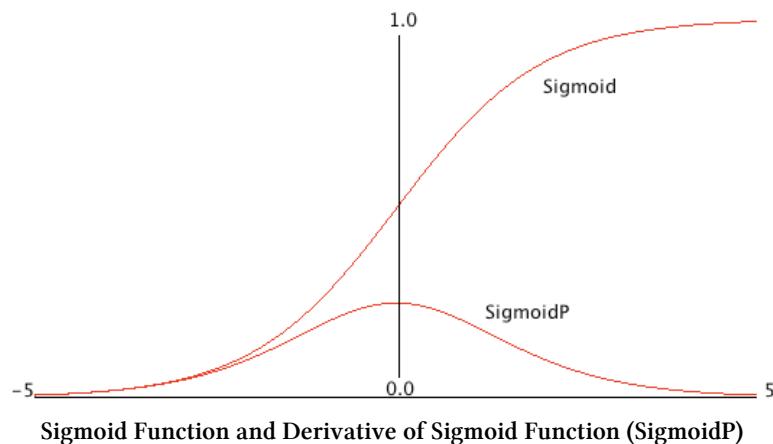
The following figure shows a simple backpropagation network with one hidden layer. Neurons in adjacent layers are connected by floating point connection strength weights. These weights start out as small random values that change as the network is trained. Weights are represented in the

following figure by arrows; in the code the weights connecting the input to the output neurons are represented as a two-dimensional array.



#### Example Backpropagation network with One Hidden Layer

Each non-input neuron has an activation value that is calculated from the activation values of connected neurons feeding into it, gated (adjusted) by the connection weights. For example, in the above figure, the value of Output 1 neuron is calculated by summing the activation of Input 1 times weight  $W_{1,1}$  and Input 2 activation times weight  $W_{2,1}$  and applying a “squashing function” like Sigmoid or Relu (see figures below) to this sum to get the final value for Output 1’s activation value. We want to flatten activation values to a relatively small range but still maintain relative values. To do this flattening we use the Sigmoid function that is seen in the next figure, along with the derivative of the Sigmoid function which we will use in the code for training a network by adjusting the weights.



Simple neural network architectures with just one or two hidden layers are easy to train using backpropagation and I have from scratch code (using no libraries) for this several of my previous books. However, here we are using Hy to write models using the TensorFlow framework which has the huge advantage that small models you experiment with on your laptop can be scaled to more parameters (usually this means more neurons in hidden layers which increases the number of weights in a model) and run in the cloud using multiple GPUs.

Except for pedantic purposes, I now never write neural network code from scratch. I take instead advantage of the many person-years of engineering work put into the development of frameworks like TensorFlow, PyTorch, mxnet, etc. We now move on to two examples built with TensorFlow.

## Deep Learning

Deep Learning models are generally understood to have many more hidden layers than simple multi-layer perceptron neural networks and often comprise multiple simple models combined together in series or in parallel. Complex architectures can be iteratively developed by manually adjusting the size of model components, changing the components, etc. Alternatively, model architecture search can be automated. At Capital One I used Google's [AdaNet project<sup>40</sup>](#) that efficiently searches for effective model architectures inside a single TensorFlow session. Now all major cloud compute provides support some form of AutoML. You need to make a decision for yourself how much effort you want to put into deeply understanding the technology, or simply learning how to use pre-trained models.

---

<sup>40</sup><https://github.com/tensorflow/adanet>

# Natural Language Processing Using Apple's Natural Language Framework

I have been working in the field of Natural Language Processing (NLP) since 1985 so I ‘lived through’ the revolutionary change in NLP that has occurred since 2014: Deep Learning results out-classed results from previous symbolic methods.

<https://developer.apple.com/documentation/naturallanguage>

I will not cover older symbolic methods of NLP here, rather I refer you to my previous books [Practical Artificial Intelligence Programming With Java<sup>41</sup>](#), [Loving Common Lisp](#), or the Savvy Programmer’s Secret Weapon<sup>42</sup>, and [Haskell Tutorial and Cookbook<sup>43</sup>](#) for examples. We get better results using Deep Learning (DL) for NLP and the libraries that Apple provides.

You will learn how to apply both DL and NLP by using the state-of-the-art full-feature libraries that Apple provides in their iOS and macOS development tools.

## Using Apple's NaturalLanguage Swift Library

We will use one of Apple’s NLP libraries consisting of pre-built models in the last chapter of this book. In order to fully understand the example in the last chapter you will need to read Apple’s high-level discussion of using CoreML [45](https://developer.apple.com/documentation/coreml<sup>44</sup></a> and their specific support for NLP <a href=).

There are many pre-trained CoreML compatible models on the web, both from Apple and also from third party (e.g., [Apple also provides tools for converting TensorFlow and PyTorch models to be compatible with CoreML \[## A simple Wrapper Library for Apple’s NLP Models\]\(https://coremltools.readme.io/docs<sup>47</sup></a>.</p></div><div data-bbox=\)](https://github.com/likedan/Awesome-CoreML-Models<sup>46</sup></a>).</p></div><div data-bbox=)

I will not go into too much detail here but I created a small wrapper library for Apple’s NLP models that will make it easier for you to jump in and have fun experimenting with them: [---

<sup>41</sup><https://leanpub.com/javaai>](https://github.com/mark-watson/Nlp_swift<sup>48</sup></a>.</p></div><div data-bbox=)

<sup>42</sup><https://leanpub.com/lovinglisp>

<sup>43</sup><https://leanpub.com/haskell-cookbook>

<sup>44</sup><https://developer.apple.com/documentation/coreml>

<sup>45</sup><https://developer.apple.com/documentation/naturallanguage/>

<sup>46</sup><https://github.com/likedan/Awesome-CoreML-Models>

<sup>47</sup><https://coremltools.readme.io/docs>

<sup>48</sup>[https://github.com/mark-watson/Nlp\\_swift](https://github.com/mark-watson/Nlp_swift)

The main library implementation file uses the `@available(OSX 10.13, *)` attribute to indicate that the following function is available on macOS 10.13 and later versions.

```

1 import Foundation
2 import NaturalLanguage
3
4 let tagger = NSLinguisticTagger(tagSchemes:[.tokenType, .language, .lexicalClass,
5 .nameType, .lemma], options: 0)
6 let options: NSLinguisticTagger.Options = [.omitPunctuation, .omitWhitespace,
7 .joinNames]
8
9 @available(OSX 10.13, *)
10 public func getEntities(for text: String) -> [(String, String)] {
11 var words: [(String, String)] = []
12 tagger.string = text
13 let range = NSRange(location: 0, length: text.utf16.count)
14 tagger.enumerateTags(in: range, unit: .word, scheme: .nameType,
15 options: options) { tag, tokenRange, stop in
16 let word = (text as NSString).substring(with: tokenRange)
17 words.append((word, tag?.rawValue ?? "unkown"))
18 }
19 return words
20 }
21
22 @available(OSX 10.13, *)
23 public func getLemmas(for text: String) -> [(String, String)] {
24 var words: [(String, String)] = []
25 tagger.string = text
26 let range = NSRange(location: 0, length: text.utf16.count)
27 tagger.enumerateTags(in: range, unit: .word, scheme: .lemma,
28 options: options) { tag, tokenRange, stop in
29 let word = (text as NSString).substring(with: tokenRange)
30 words.append((word, tag?.rawValue ?? "unkown"))
31 }
32 return words
33 }
```

The public function `getEntities` takes a `String` parameter called `text` and returns an array of tuples containing `(String, String)`. Here's a breakdown of what this function does:

- The function initializes an empty array called `words` to store the extracted entities.
- The line `tagger.string = text` sets the input text for a `tagger` object. The `tagger` is an instance of `NSLinguisticTagger`, which is a natural language processing class provided by Apple's Foundation framework.

- The next line creates an **NSRange** object called **range** that represents the entire length of the input text.
- The **tagger.enumerateTags(in:range, unit:.word, scheme:.nameType, options:options)** method is called to iterate over the words in the input text and extract their associated tags. The **in:** parameter specifies the range of the text to process. The **unit:** parameter specifies that the enumeration should be done on a word-by-word basis. The **scheme:** parameter specifies the linguistic scheme to use, in this case, the **.nameType** scheme, which is used to identify named entities. The **options:** parameter specifies additional options or settings for the tagger.
- Inside the enumeration block, the code retrieves the current word and its associated tag using the **tokenRange** and **tag** parameters.
- The line **let word = (text as NSString).substring(with: tokenRange)** extracts the substring corresponding to the current word using **tokenRange**.
- The line **words.append((word, tag?.rawValue ?? "unknown"))** appends a tuple containing the extracted word and its associated tag to the **words** array. If the tag is nil, it uses the default value of “unknown”.
- Finally, the **words** array is returned, which contains all the extracted entities (words and their associated tags) from the input text.

The public function called **getLemmas** that takes a String parameter called **text** and returns an array of tuples containing **(String, String)**. Here's a breakdown of what the function **getLemmas** is very similar to the last function **getEntities**. The function **getLemmas** does the following:

- The function initializes an empty array called **words** to store the extracted lemmas.
- The line **tagger.string = text** sets the input text for a tagger object.
- The next line creates an **NSRange** object called **range** that represents the entire length of the input text.
- The **tagger.enumerateTags(in:range, unit:.word, scheme:.lemma, options:options)** method is called to iterate over the words in the input text and extract their corresponding lemmas.
- Inside the enumeration block, the code retrieves the current word and its associated lemma using the **tokenRange** and **tag** parameters.
- The line **let word = (text as NSString).substring(with: tokenRange)** extracts the substring corresponding to the current word using **tokenRange**.
- Finally, the **words** array is returned, which contains all the extracted lemmas (words and their associated base forms) from the input text.

In summary, function **getLemmas** uses the **NSLinguisticTagger** to perform linguistic analysis on a given text and extract the base forms (lemmas) of words. The lemmas are then stored in an array of tuples and returned as the result of the function.

Here is some test code:

```
1 let quote = "President George Bush went to Mexico with IBM representatives. Here's t\
2 o the crazy ones. The misfits. The rebels. The troublemakers. The round pegs in the \
3 square holes. The ones who see things differently. They're not fond of rules. And th\
4 ey have no respect for the status quo. You can quote them, disagree with them, glori\
5 fy or vilify them. About the only thing you can't do is ignore them. Because they ch\
6 ange things. They push the human race forward. And while some may see them as the cr\
7 azy ones, we see genius. Because the people who are crazy enough to think they can c\
8 hange the world, are the ones who do. - Steve Jobs (Founder of Apple Inc.)"
9 if #available(OSX 10.13, *) {
10 print("\nEntities:\n")
11 print(getEntities(for: quote))
12 print("\nLemmas:\n")
13 print(getLemmas(for: quote))
14 }
```

Here is an edited listing of the output with most of the output removed for brevity:

```
1 Entities:
2
3 [("President", "OtherWord"), ("George Bush", "PersonalName"), ("went", "OtherWord"), \
4 ("to", "OtherWord"), ("Mexico", "PlaceName"), ("with", "OtherWord"), ("IBM", "Organ\
5 izationName"),
6 ...
7
8 Lemmas:
9
10 [("President", "President"), ("George Bush", "George"), ("went", "go"), ("to", "to")\
11 , ("Mexico", "Mexico"),
12 ...]
```

# Documents Question Answering Using OpenAI GPT4 APIs and a Local Embeddings Vector Database

The examples in this chapter are inspired by the Python LangChain and LlamaIndex projects, with just the parts I need for my projects written from scratch in Common Lisp. I wrote a Python book “LangChain and LlamaIndex Projects Lab Book: Hooking Large Language Models Up to the Real World Using GPT-3, ChatGPT, and Hugging Face Models in Applications” in March 2023: <https://leanpub.com/langchain> that you might also be interested in.

The GitHub repository for this example can be found here: [https://github.com/mark-watson/Docs\\_QA\\_Swift](https://github.com/mark-watson/Docs_QA_Swift)<sup>49</sup>.

The entire example is in one Swift source file `main.swift`. All of the program listings in this chapter can be found in this single source file.

We use two models in this example: a vector embedding model and a gpt-4o-mini conversation model (see bottom of this file). The vector embedding model is used to generate a vector embedding. The gpt-4o-mini model is used to generate a response to a prompt. The vector embedding model is used to compare the similarity of two prompts.

## Extending the String Class

```
1 import Foundation
2 import NaturalLanguage
3
4 // String utilities:
5
6 extension String {
7 func removeCharacters(from forbiddenChars: CharacterSet) -> String {
8 let passed = self.unicodeScalars.filter { !forbiddenChars.contains($0) }
9 return String(String.UnicodeScalarView(passed))
10 }
11
12 func removeCharacters(from: String) -> String {
13 return removeCharacters(from: CharacterSet(charactersIn: from))
```

---

<sup>49</sup>[https://github.com/mark-watson/Docs\\_QA\\_Swift](https://github.com/mark-watson/Docs_QA_Swift)

```

14 }
15 func plainText() -> String {
16 return self.removeCharacters(from:
17 "\u00a0()%$#@[]{}<>").replacingOccurrences(of: "\n\
18 ",
19 with: " "))
20 }
21 }
```

## Implementing a Local Vector Database for Document Embeddings

```

1 let openai_key = ProcessInfo.processInfo.environment["OPENAI_KEY"]!
2
3 let openAiHost = "https://api.openai.com/v1/embeddings"
4
5 func openAiHelper(body: String) -> String {
6 var ret = ""
7 var content = "{}"
8 let requestUrl = URL(string: openAiHost)!
9 var request = URLRequest(url: requestUrl)
10 request.httpMethod = "POST"
11 request.httpBody = body.data(using: String.Encoding.utf8);
12 request.setValue("application/json", forHTTPHeaderField: "Content-Type")
13 request.setValue("Bearer " + openai_key, forHTTPHeaderField: "Authorization")
14 let task = URLSession.shared.dataTask(with: request) { (data, response, error) in
15 if let error = error {
16 print("--> Error accessing OpenAI servers: \(error)")
17 return
18 }
19 if let data = data, let s = String(data: data, encoding: .utf8) {
20 content = s
21 //print("** s=\", s")
22 CFRunLoopStop(CFRunLoopGetMain())
23 }
24 }
25 task.resume()
26 CFRunLoopRun()
27 let c = String(content)
28 let i1 = c.range(of: "\"embedding\":")
29 if let r1 = i1 {
```

```
30 let i2 = c.range(of: "]")
31 if let r2 = i2 {
32 ret = String(String(String(c[r1.lowerBound..
```

The source file contains example code for creating embeddings and using dot product work to find semantic similarity:

```
1 let emb1 = embeddings(someText: "John bought a new car")
2 let emb2 = embeddings(someText: "Sally drove to the store")
3 let emb3 = embeddings(someText: "The dog saw a cat")
4 let dotProductResult1 = dotProduct(emb1, emb2)
5 print(dotProductResult1)
6 let dotProductResult2 = dotProduct(emb1, emb3)
7 print(dotProductResult2)
```

The output is:

```
1 0.8416926
2 0.79411536
```

For this example, we use an in-memory store of embedding vectors and chunk text. A text document is broken into smaller chunks of text. Each chunk is embedded and stored in the embeddingsStore. The chunk text is stored in the chunks array. The embeddingsStore and chunks array are used to find the most similar chunk to a prompt. The most similar chunk is used to generate a response to the prompt.

```
1 var embeddingsStore: Array<[Float]> = Array()
2 var chunks: Array<String> = Array()
3
4 func addEmbedding(_ embedding: [Float]) {
5 embeddingsStore.append(embedding)
6 //print("Added embedding: count=\(embeddingsStore.count) \nembedding")
7 }
8
9 func addChunk(_ chunk: String) {
10 chunks.append(chunk)
11 }
```

## Create Local Embeddings Vectors From Local Text Files With OpenAI GPT APIs

```
1 func readList(_ input: String) -> [Float] {
2 return input.split(separator: ",\n").compactMap {
3 Float($0.trimmingCharacters(in: .whitespaces))
4 }
5 }
6
7 let fileManager = FileManager.default
8 let currentDirectoryURL = URL(fileURLWithPath: fileManager.currentDirectoryPath)
9 let dataDirectoryURL = currentDirectoryURL.appendingPathComponent("data")
10
11 // Top level code expression to process all *.txt files in the data/ directory:
12
13 do {
14 let directoryContents = try fileManager.contentsOfDirectory(at: dataDirectoryURL \
15 , includingPropertiesForKeys: nil)
16 let txtFiles = directoryContents.filter { $0.pathExtension == "txt" }
17 for txtFile in txtFiles {
18 let content = try String(contentsOf: txtFile)
19 let chnks = segmentTextIntoChunks(text: content/plainText(),
20 max_chunk_size: 100)
21 for chunk in chnks {
22 let embedding = embeddings(someText: chunk)
23 if embedding.count > 0 {
24 addEmbedding(embedding)
25 addChunk(chunk)
26 }
27 }
28 }
29 } catch {
30 }
31
32 func segmentTextIntoSentences(text: String) -> [String] {
33 let tokenizer = NLTokenizer(unit: .sentence)
34 tokenizer.string = text
35 let sentences = tokenizer.tokens(for: text.startIndex..
```

```

44 var chunks: Array<String> = Array()
45 var currentChunk = ""
46 var currentChunkSize = 0
47 for sentence in sentences {
48 if currentChunkSize + sentence.count < max_chunk_size {
49 currentChunk += sentence
50 currentChunkSize += sentence.count
51 } else {
52 chunks.append(currentChunk)
53 currentChunk = sentence
54 currentChunkSize = sentence.count
55 }
56 }
57 return chunks
58 }
```

## Using Local Embeddings Vector Database With OpenAI GPT APIs

We use the OpenAI QA API using gpt-4o-mini model (reformatted to fit the page width):

```

1 let openAiQaHost = "https://api.openai.com/v1/chat/completions"
2
3 func openAiQaHelper(body: String) -> String {
4 var ret = ""
5 var content = "{}"
6 let requestUrl = URL(string: openAiQaHost)!
7 var request = URLRequest(url: requestUrl)
8 request.httpMethod = "POST"
9 request.httpBody = body.data(using: String.Encoding.utf8);
10 request.setValue("application/json", forHTTPHeaderField: "Content-Type")
11 request.setValue("Bearer " + openai_key, forHTTPHeaderField: "Authorization")
12 let task = URLSession.shared.dataTask(with: request) { (data, response, error) in
13 if let error = error {
14 print("--> Error accessing OpenAI servers: \(error)")
15 return
16 }
17 if let data = data, let s = String(data: data, encoding: .utf8) {
18 content = s
19 CFRunLoopStop(CFRunLoopGetMain())
```

```

20 }
21 }
22 task.resume()
23 CFRunLoopRun()
24 let c = String(content)
25 //print("DEBUG response c:", c)
26 // pull returned content for string instead of using a
27 // JSON parser:
28 let i1 = c.range(of: "\"content\":")
29 if let r1 = i1 {
30 let i2 = c.range(of: "\"}")
31 if let r2 = i2 {
32 ret = String(
33 String(
34 String(c[r1.lowerBound..

```

```

63 let dotProductResult = dotProduct(queryEmbedding, embeddingsStore[i])
64 if dotProductResult > 0.8 {
65 contextText.append(chunks[i])
66 contextText.append(" ")
67 }
68 }
69 //print("\n\n++++++ contextText = \\" + contextText + "\\n\\n")
70 let answer = questionAnswering(context: contextText, question: query)
71 //print("* * debug: query: ", query)
72 //print("* * debug: answer:", answer)
73 return answer
74
75 }
76
77 print(query("What is the history of chemistry?"))
78 print(query("What is the definition of sports?"))
79 print(query("What is the microeconomics?"))

```

The output for these three questions looks like:

1 The history of chemistry dates back to ancient times when people began to manipulate\\  
 2 materials to produce useful products. The ancient Egyptians were skilled in metallu\\  
 3 rgy and used various chemicals to embalm bodies. The Greeks were interested in theor\\  
 4 ies of matter and sought to understand the nature of substances.\\n\\nDuring the Middl\\  
 5 e Ages, alchemy became popular, with alchemists seeking to transform base metals int\\  
 6 o gold and searching for an elixir of life. While alchemy was considered a pseudosci\\  
 7 ence, it did lead to important discoveries such as the distillation of alcohol and t\\  
 8 he discovery of various acids.\\n\\nThe Scientific Revolution of the 17th century brou\\  
 9 ght about significant changes in chemistry. The work of Robert Boyle, Antoine Lavois\\  
 10 ier, and others laid the foundation for modern chemistry. Lavoisier is considered th\\  
 11 e father of modern chemistry for his work in establishing the law of conservation of\\  
 12 mass, which states that matter cannot be created or destroyed.\\n\\nThe 19th century \\  
 13 saw the development of organic chemistry, as scientists sought to understand the che\\  
 14 mistry of carbon-based compounds, which make up many biological molecules. The 20th \\  
 15 century brought about significant advances in analytical chemistry, as well as the d\\  
 16 evelopment of quantum mechanics and the discovery of the structure of DNA, which rev\\  
 17 olutionized the field of biochemistry.\\n\\nToday, chemistry plays a critical role in \\  
 18 fields such as medicine, agriculture, materials science, and environmental science.  
 19  
 20  
 21 Sports can be defined as activities involving physical athleticism, physical dexteri\\  
 22 ty, and governed by rules to ensure fair competition and consistent adjudication of \\  
 23 the winner. The term \"sport\" originally meant leisure, but it now primarily refers\\

24 to physical activities that involve competition at various levels of skill and prof\\  
25 ciency. Some organizations also include all physical activity and exercise in the d\\  
26 efinition of sport.

27

28

29 Microeconomics is a branch of economics that focuses on the behavior and decision-ma\\  
30 king of individual units within an economy, such as households, firms, and industrie\\  
31 s. It examines how these units interact in various markets to determine the prices o\\  
32 f goods and services and how resources are allocated efficiently. Microeconomics als\\  
33 o considers the role of government policies and regulations in influencing these int\\  
34 eractions and outcomes. Topics studied in microeconomics include supply and demand, \\  
35 market structures, consumer behavior, production and cost analysis, and welfare anal\\  
36 ysis.

## Wrap Up for Using Local Embeddings Vector Database to Enhance the Use of GPT3 APIs With Local Documents

As I write this in early April 2023, I have been working almost exclusively with OpenAI APIs for the last year and using the Python libraries for LangChain and LlamaIndex for the last three months.

I started writing the examples in this chapter for my own use, implementing a tiny subset of the LangChain and LlamaIndex libraries in Swift in order to write efficient command line utilities for creating local embedding vector data stores and for interactive chat using my own data.

By writing about my “scratching my own itch” command line experiments here I hope that I get pull requests for [https://github.com/mark-watson/Docs\\_QA\\_Swift](https://github.com/mark-watson/Docs_QA_Swift) from readers who are interested in helping to extend this code with new functionality.

# **Part 4: Knowledge Representation and Data Acquisition**

In this part we cover:

- Introduction to the semantic web and linked data
- A general discussion of Knowledge Representation
- Create Knowledge Graphs from text input
- Knowledge Graph Explorer application

# Linked Data and the Semantic Web

Tim Berners Lee, James Hendler, and Ora Lassila wrote in 2001 an article for Scientific American where they introduced the term Semantic Web. Here I do not capitalize semantic web and use the similar term linked data somewhat interchangeably with semantic web.

In the same way that the web allows links between related web pages, linked data supports linking associated data on the web together. I view linked data as a relatively simple way to specify relationships between data sources on the web while the semantic web has a much larger vision: the semantic web has the potential to be the entirety of human knowledge represented as data on the web in a form that software agents can work with to answer questions, perform research, and to infer new data from existing data.

While the “web” describes information for human readers, the semantic web is meant to provide structured data for ingestion by software agents. This distinction will be clear as we compare WikiPedia, made for human readers, with DBPedia which uses the info boxes on WikiPedia topics to automatically extract RDF data describing WikiPedia topics. Let’s look at the WikiPedia topic for the town I live in Sedona, Arizona, and show how the info box on the English version of the [WikiPedia topic page for Sedona](https://en.wikipedia.org/wiki/Sedona,_Arizona) [https://en.wikipedia.org/wiki/Sedona,\\_Arizona<sup>50</sup>](https://en.wikipedia.org/wiki/Sedona,_Arizona) maps to the [DBPedia page](http://dbpedia.org/page/Sedona,_Arizona) [http://dbpedia.org/page/Sedona,\\_Arizona<sup>51</sup>](http://dbpedia.org/page/Sedona,_Arizona). Please open both of these WikiPedia and DBPedia URIs in two browser tabs and keep them open for reference.

I assume that the format of the WikiPedia page is familiar so let’s look at the DBPedia page for Sedona that in human readable form shows the RDF statements with Sedona Arizona as the subject. RDF is used to model and represent data. RDF is defined by three values so an instance of an RDF statement is called a *triple* with three parts:

- subject: a URI (also referred to as a “Resource”)
- property: a URI (also referred to as a “Resource”)
- value: a URI (also referred to as a “Resource”) or a literal value (like a string or a number with optional units)

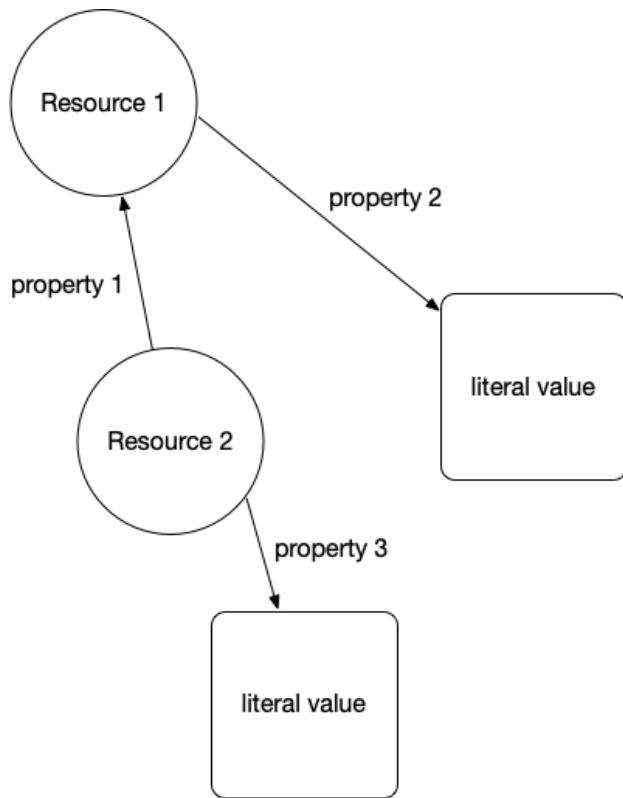
The subject for each Sedona related triple is the above URI for the DBPedia human readable page. The subject and property references in an RDF triple will almost always be a URI that can ground an entity to information on the web. The human readable page for Sedona lists several properties and the values of these properties. One of the properties is “dbo:areaCode” where “dbo” is a name space reference (in this case for a [DatatypeProperty<sup>52</sup>](#)).

<sup>50</sup>[https://en.wikipedia.org/wiki/Sedona,\\_Arizona](https://en.wikipedia.org/wiki/Sedona,_Arizona)

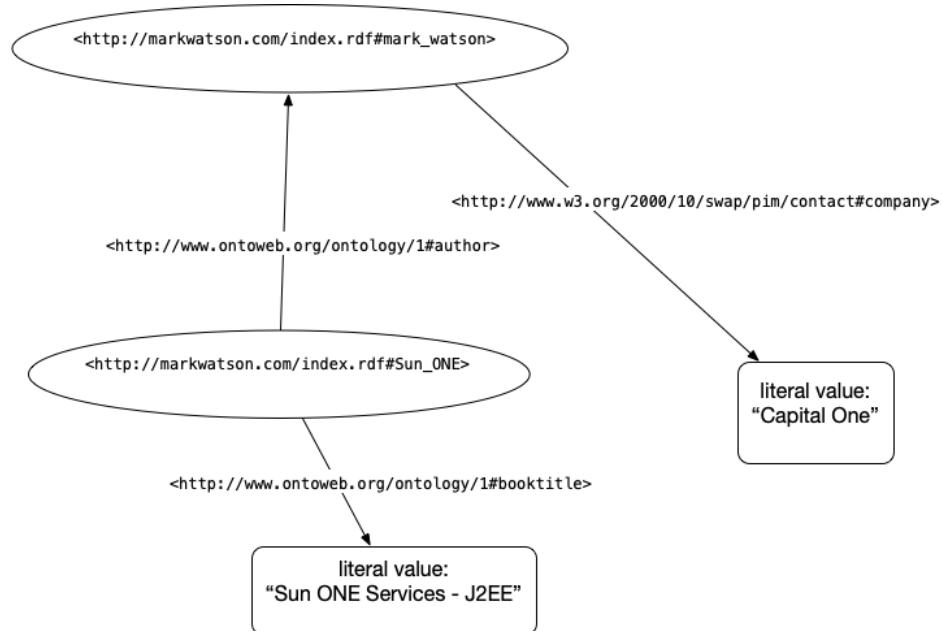
<sup>51</sup>[http://dbpedia.org/page/Sedona,\\_Arizona](http://dbpedia.org/page/Sedona,_Arizona)

<sup>52</sup><http://www.w3.org/2002/07/owl#DatatypeProperty>

The following two figures show an abstract representation of linked data and then a sample of linked data with actual web URIs for resources and properties:



Abstract RDF representation with 2 Resources, 2 literal values, and 3 Properties



Concrete example using RDF seen in last chapter showing the RDF representation with 2 Resources, 2 literal values, and 3 Properties

We will use the SPARQL query language (SPARQL for RDF data is similar to SQL for relational database queries). Let's look at an example using the RDF in the last figure:

```

1 "select ?v where { <http://markwatson.com/index.rdf#Sun_ONE>
2 <http://www.ontoweb.org/ontology/1#booktitle>
3 ?v }

```

This query should return the result “Sun ONE Services - J2EE”. If you wanted to query for all URI resources that are books with the literal value of their titles, then you can use:

```

1 "select ?s ?v where { ?s
2 <http://www.ontoweb.org/ontology/1#booktitle>
3 ?v }

```

Note that **?s** and **?v** are arbitrary query variable names, here standing for “subject” and “value”. You can use more descriptive variable names like:

```

1 "select ?bookURI ?bookTitle where
2 { ?bookURI
3 <http://www.ontoweb.org/ontology/1#booktitle>
4 ?bookTitle }

```

We will be diving a little deeper into RDF examples in the next chapter when we write a tool for using RDF data from DBpedia to find information about entities (e.g., people, places, organizations)

and the relationships between entities. For now I want you to understand the idea of RDF statements represented as triples, that web URIs represent things, properties, and sometimes values, and that URIs can be followed manually (often called “dereferencing”) to see what they reference in human readable form.

## Understanding the Resource Description Framework (RDF)

Text data on the web has some structure in the form of HTML elements like headers, page titles, anchor links, etc. but this structure is too imprecise for general use by software agents. RDF is a method for encoding structured data in a more precise way.

RDF specifies graph structures and can be serialized for storage or for service calls in XML, Turtle, N3, and other formats. I like the Turtle format and suggest that you pause reading this book for a few minutes and look at this World Wide Web Consortium Turtle RDF primer at <https://www.w3.org/2007/02/turtle/primer/><sup>53</sup>.

## Frequently Used Resource Namespaces

The following standard namespaces are frequently used:

- RDF <https://www.w3.org/TR/rdf-syntax-grammar/><sup>54</sup>
- RDFS <https://www.w3.org/TR/rdf-schema/><sup>55</sup>
- OWL <http://www.w3.org/2002/07/owl#><sup>56</sup>
- XSD <http://www.w3.org/2001/XMLSchema#><sup>57</sup>
- FOAF <http://xmlns.com/foaf/0.1/><sup>58</sup>
- SKOS <http://www.w3.org/2004/02/skos/core#><sup>59</sup>
- DOAP <http://usefulinc.com/ns/doap#><sup>60</sup>
- DC <http://purl.org/dc/elements/1.1/><sup>61</sup>
- DCTERMS <http://purl.org/dc/terms/><sup>62</sup>
- VOID <http://rdfs.org/ns/void#><sup>63</sup>

Let's look into the Friend of a Friend (FOAF) namespace. Click on the above link for FOAF <http://xmlns.com/foaf/0.1/><sup>64</sup> and find the definitions for the FOAF Core:

---

<sup>53</sup><https://www.w3.org/2007/02/turtle/primer/>

<sup>54</sup><https://www.w3.org/TR/rdf-syntax-grammar/>

<sup>55</sup><https://www.w3.org/TR/rdf-schema/>

<sup>56</sup><http://www.w3.org/2002/07/owl#>

<sup>57</sup><http://www.w3.org/2001/XMLSchema#>

<sup>58</sup><http://xmlns.com/foaf/0.1/>

<sup>59</sup><http://www.w3.org/2004/02/skos/core#>

<sup>60</sup><http://usefulinc.com/ns/doap#>

<sup>61</sup><http://purl.org/dc/elements/1.1/>

<sup>62</sup><http://purl.org/dc/terms/>

<sup>63</sup><http://rdfs.org/ns/void#>

<sup>64</sup><http://xmlns.com/foaf/0.1/>

```
1 Agent
2 Person
3 name
4 title
5 img
6 depiction (depicts)
7 familyName
8 givenName
9 knows
10 based_near
11 age
12 made (maker)
13 primaryTopic (primaryTopicOf)
14 Project
15 Organization
16 Group
17 member
18 Document
19 Image
```

and for the Social Web:

```
1 mbox
2 homepage
3 weblog
4 openid
5 jabberID
6 mbox_sha1sum
7 interest
8 topic_interest
9 topic (page)
10 workplaceHomepage
11 workInfoHomepage
12 schoolHomepage
13 publications
14 currentProject
15 pastProject
16 account
17 OnlineAccount
18 accountName
19 accountServiceHomepage
20 PersonalProfileDocument
21 tipjar
```

```

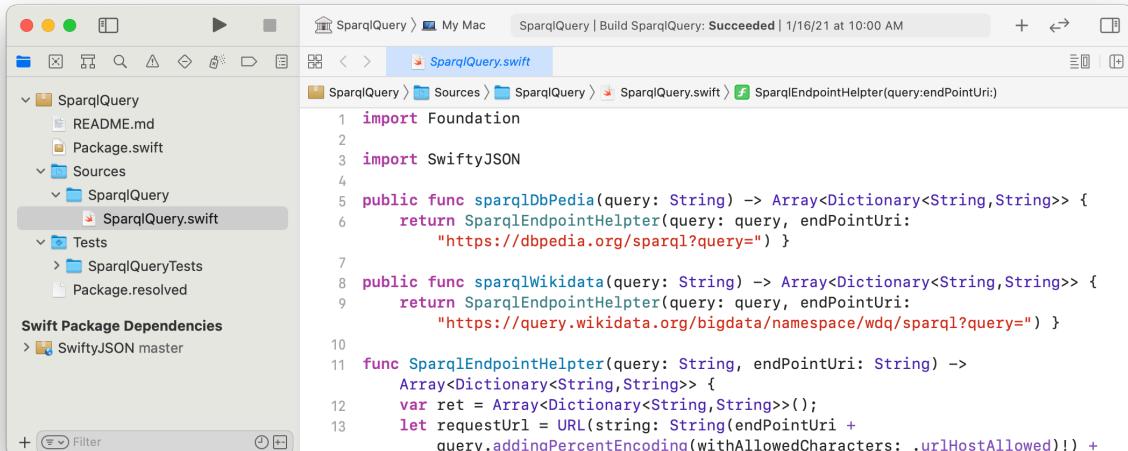
22 sha1
23 thumbnail
24 logo

```

You now have seen a few common Schemas for RDF data. Another Schema that is widely used for annotating web sites that we won't need for our examples here, is [schema.org](https://schema.org)<sup>65</sup>.

## Understanding the SPARQL Query Language

For the purposes of the material in this book, the two sample SPARQL queries here are sufficient for you to get started using my SPARQL library [https://github.com/mark-watson/SparqlQuery\\_swift](https://github.com/mark-watson/SparqlQuery_swift)<sup>66</sup> with arbitrary RDF data sources and simple queries.



My Swift SPARQL library open in Xcode

The Apache Foundation has a [good introduction to SPARQL](#)<sup>67</sup> that I refer you to for more information.

## Semantic Web and Linked Data Wrap Up

In the next chapter we will use natural language processing to extract structured information from raw text from SPARQL queries. We will be using my Swift SPARQL library [https://github.com/mark-watson/SparqlQuery\\_swift](https://github.com/mark-watson/SparqlQuery_swift)<sup>68</sup> as well as two pre-trained CoreML deep learning models.

<sup>65</sup><https://schema.org>

<sup>66</sup>[https://github.com/mark-watson/SparqlQuery\\_swift](https://github.com/mark-watson/SparqlQuery_swift)

<sup>67</sup><https://jena.apache.org/tutorials/sparql.html>

<sup>68</sup>[https://github.com/mark-watson/SparqlQuery\\_swift](https://github.com/mark-watson/SparqlQuery_swift)

# Example Application: iOS and macOS Versions of my KnowledgeBookNavigator

I used many of the techniques discussed in this book, the Swift language, and the SwiftUI user interface framework to develop Swift version of my Knowledge Graph Navigator application for macOS. I originally wrote this as an example program in Common Lisp for another book project.

The GitHub repository for the KGN example is <https://github.com/mark-watson/KGN><sup>69</sup>. I copied the code from my stand-alone Swift libraries to this example to make it self contained. The easiest way to browse the source code is to open this project in Xcode.

I submitted the KGN app that we discuss in this chapter to Apple's store and is available as a macOS app. If you load this project into Xcode, you can also build and run the iOS and iPadOS targets.

You will need to have read through the last chapter on semantic web and linked data technologies to understand this example because quite a lot of the code has embedded SPARQL queries to get information from [DBpedia.org](https://dbpedia.org)<sup>70</sup>.

The other major part of this app is a slightly modified version of Apple's question answering (QA) example using the BERT model in CoreML. Apple's code is in the subdirectory **AppleBERT**. Please read the README file for this project and follow the directions for downloading and using Apple's model and vocabulary file.

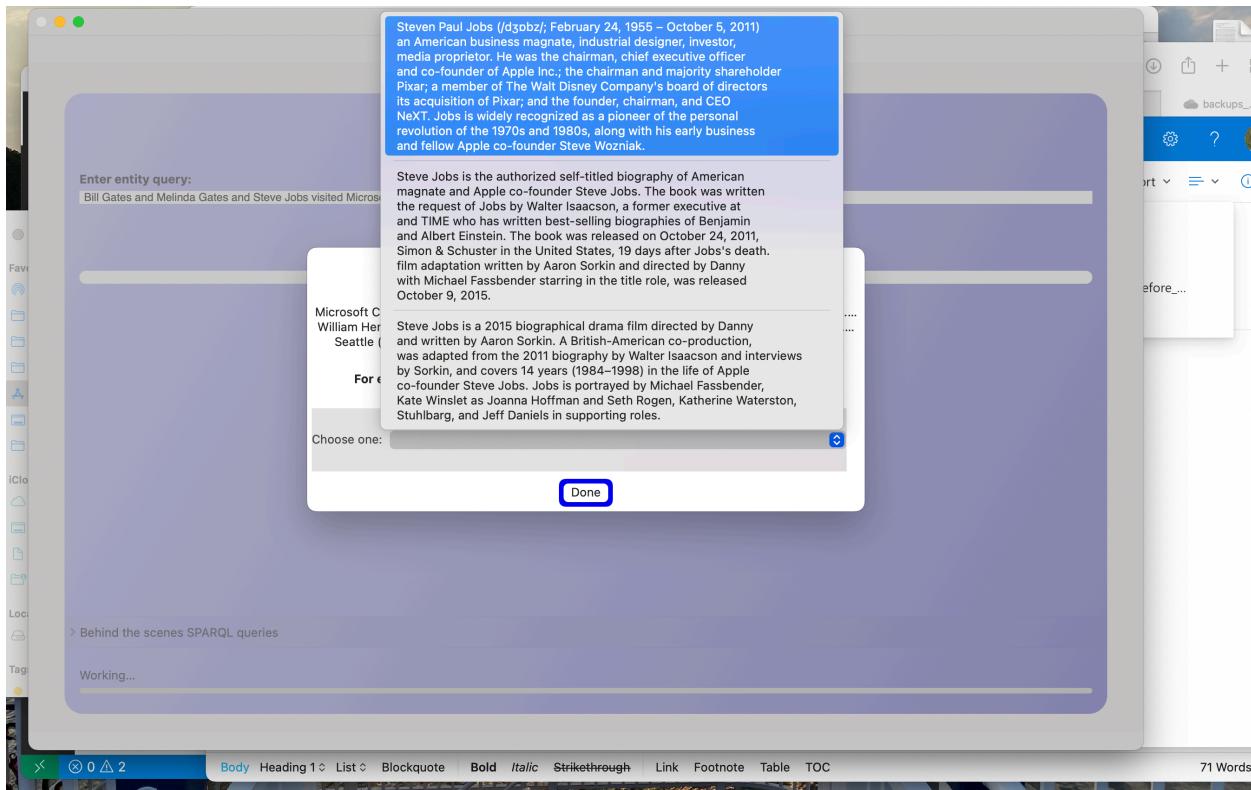
## Screen Shots of macOS Application

In the first screenshot seen below, I had entered query text that included "Steve Jobs" and the popup list selector is used to let the user select which "Steve Jobs" entity from DBpedia that they want to use.

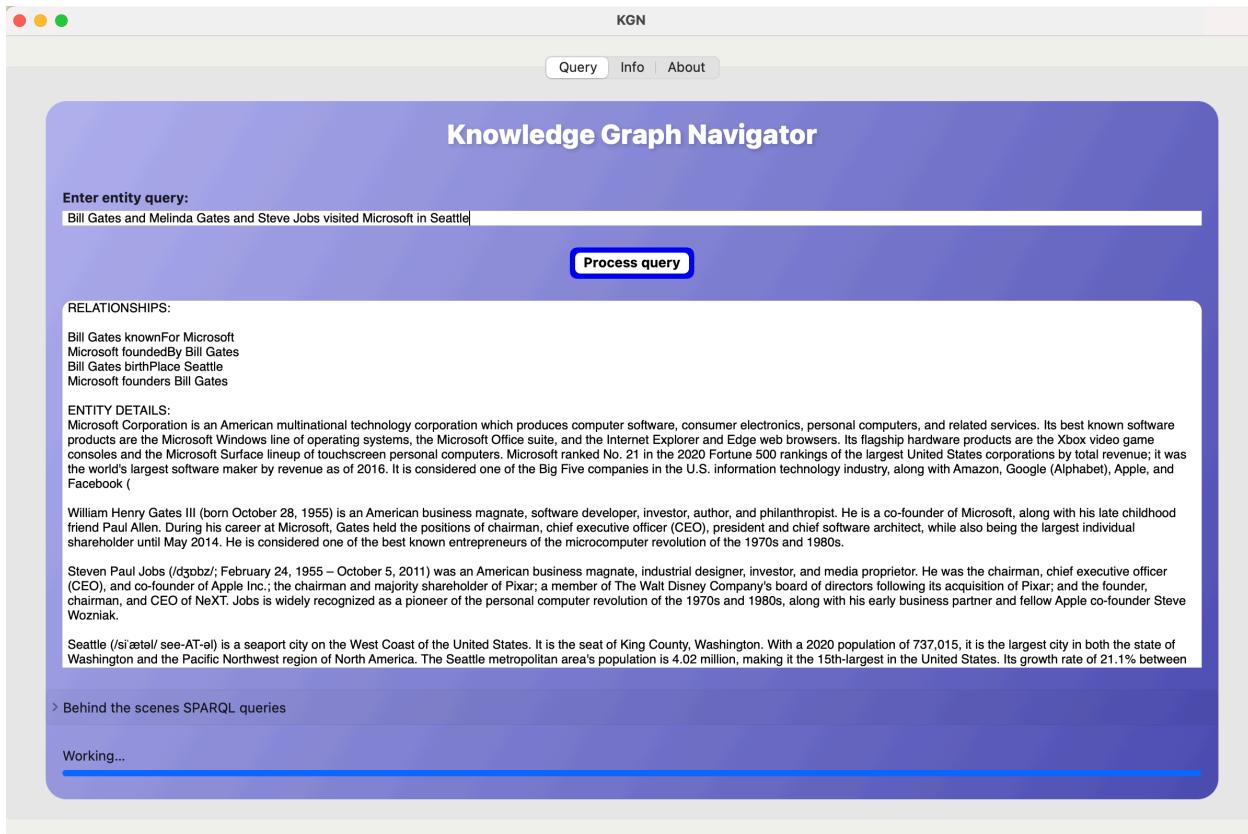
---

<sup>69</sup><https://github.com/mark-watson/KGN>

<sup>70</sup><https://dbpedia.org>



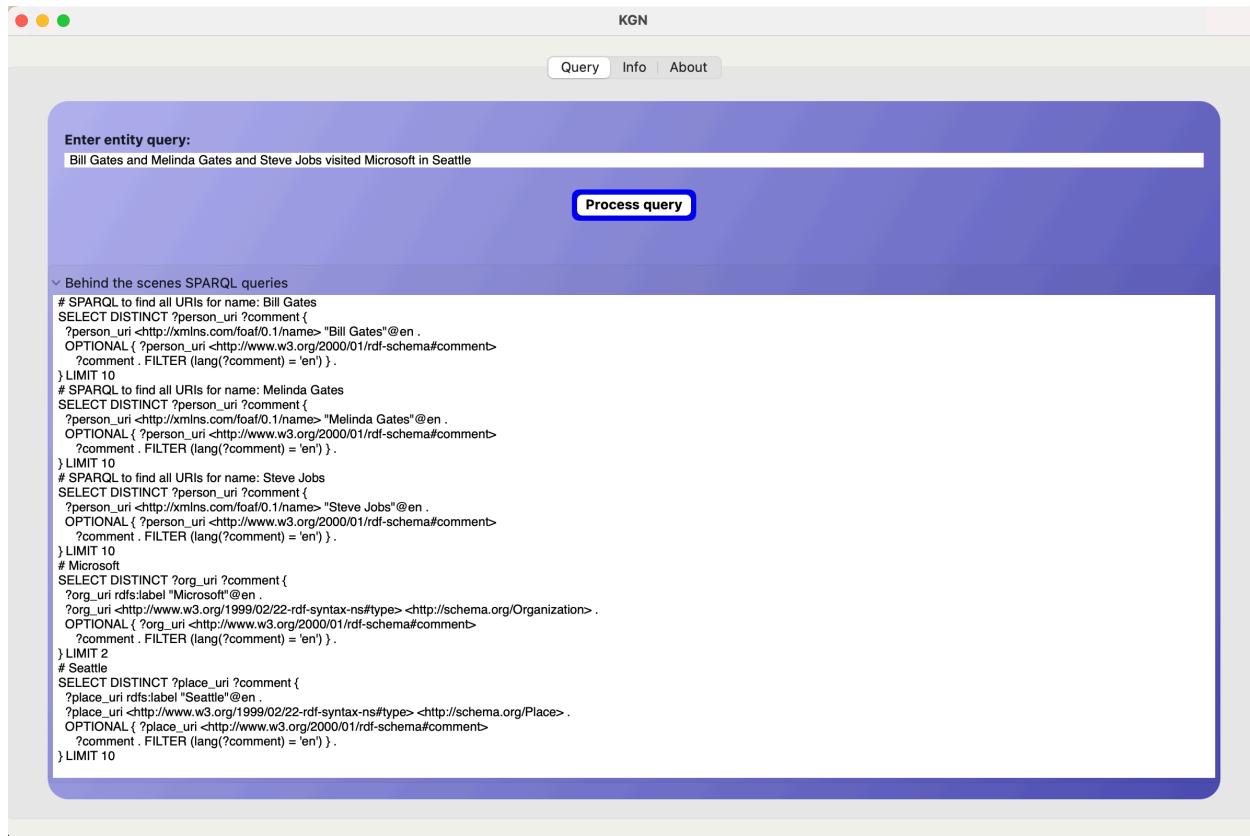
Entered query and KGN is asking user to disambiguate which “Steve Jobs” they want information for



### Showing results

The previous screenshot shows the results to the query displayed as English text.

Notice the app prompt “Behind the scenes SPARQL queries” near the bottom of the app window. If you click on this field then the SPARQL queries used to answer the question are shown, as on the next screenshot:



Showing SPARQL queries used to gather data

## Application Code Listings

I will list some of the code for this example application and I suggest that you, dear reader, also open this project in Xcode in order to navigate the sample code and more carefully read through it.

### SPARQL

I introduced you to the use of SPARQL in the last chapter. This library can be used by adding a reference to the `Project.swift` file for this project. You can also clone the [GitHub repository](#) [`https://github.com/mark-watson/Nlp\_swift`](https://github.com/mark-watson/Nlp_swift)<sup>71</sup> to have the source code for local viewing and modification and I have copied the code into the KGN project.

The file `SparqlQuery.swift` is shown here:

<sup>71</sup>[https://github.com/mark-watson/Nlp\\_swift](https://github.com/mark-watson/Nlp_swift)

```

1 import Foundation
2
3 public func sparqlDbPedia(query: String) -> Array<Dictionary<String, String>> {
4 return SparqlEndpointHelper(query: query,
5 endPointUri: "https://dbpedia.org/sparql?query=")
6
7 public func sparqlWikidata(query: String) -> Array<Dictionary<String, String>> {
8 return SparqlEndpointHelper(query: query,
9 endPointUri:
10 "https://query.wikidata.org/bigdata/namespace/wdq/sparql?query=")
11
12 public func SparqlEndpointHelper(query: String,
13 endPointUri: String) ->
14 Array<Dictionary<String, String>> {
15 var ret = Set<Dictionary<String, String>>()
16 var content = "{}"
17
18 let maybeString = cacheLookupQuery7(key: query)
19 if maybeString?.count ?? 0 > 0 {
20 content = maybeString ?? ""
21 } else {
22 let requestUrl = URL(string: String(endPointUri + query.addingPercentEncoding\
23 g(withAllowedCharacters:
24 .urlHostAllowed)!!) + "&format=json")!
25 do { content = try String(contentsOf: requestUrl) }
26 catch let error { print(error) }
27 }
28 let json = try? JSONSerialization.jsonObject(with: Data(content.utf8),
29 options: [])
30 if let json2 = json as! Optional<Dictionary<String, Any?>> {
31 if let head = json2["head"] as? Dictionary<String, Any> {
32 if let xvars = head["vars"] as! NSArray? {
33 if let results = json2["results"] as? Dictionary<String, Any> {
34 if let bindings = results["bindings"] as! NSArray? {
35 if bindings.count > 0 {
36 for i in 0...(bindings.count-1) {
37 if let first_binding =
38 bindings[i] as? Dictionary<String,
39 Dictionary<String, String>> {
40 var ret2 = Dictionary<String, String>();
41 for key in xvars {
42 let key2 : String = key as! String
43 if let vals = (first_binding[key2]) {

```

```

44 let vv : String = vals["value"] ?? "err2"
45 ret2[key2] = vv } }
46 if ret2.count > 0 {
47 ret.insert(ret2)
48 }}}}}}}}}
49 return Array(ret) }
```

The file `QueryCache.swift` contains code written by Khoa Pham (MIT License) that can be found in the GitHub repository <https://github.com/onmyway133/EasyStash><sup>72</sup>. This file is used to cache SPARQL queries and the results. In testing this application I noticed that there were many repeated queries to DBpedia so I decided to cache results. Here is the simple API I added on top of Khoa Pham's code:

```

1 // Created by khoa on 27/05/2019.
2 // Copyright © 2019 Khoa Pham. All rights reserved. MIT License.
3 // https://github.com/onmyway133/EasyStash
4 //
5
6 import Foundation
7
8 // Mark's simple wrapper:
9
10 var storage: Storage? = nil
11
12 public func cacheStoreQuery(key: String, value: String) {
13 do { try storage?.save(object: value, forKey: key) } catch {} }
14 }
15 public func cacheLookupQuery7(key: String) -> String? {
16 // optional DEBUG code: clear cache
17 //do { try storage?.removeAll() } catch { print("ERROR CLEARING CACHE") }
18 do {
19 return try storage?.load(forKey: key, as: String.self)
20 } catch { return "" }
21 }
22
23 // remaining code not shown for brevity.
```

The code in file `GenerateSparql.swift` is used to generate queries for DBpedia. The line-wrapping for embedded SPARQL queries in the next code section is difficult to read so you may want to open the source file in Xcode. Please note that the KGN application prints out the SPARQL queries used to fetch information from DBpedia. The embedded SPARQL query templates used here have variable slots that filled in at runtime to customize the queries.

---

<sup>72</sup><https://github.com/onmyway133/EasyStash>

```

1 //
2 // GenerateSparql.swift
3 // KGNbeta1
4 //
5 // Created by Mark Watson on 2/28/20.
6 // Copyright © 2021 Mark Watson. All rights reserved.
7 //
8
9 import Foundation
10
11 public func uri_to_display_text(uri: String)
12 -> String {
13 return uri.replacingOccurrences(of: "http://dbpedia.org/resource/Category/",
14 with: "").
15 replacingOccurrences(of: "http://dbpedia.org/resource/",
16 with: "").
17 replacingOccurrences(of: "_", with: " ")
18 }
19
20 public func get_SPARQL_for_finding_URIIs_for_PERSON_NAME(nameString: String)
21 -> String {
22 return
23 "# SPARQL to find all URIs for name: " +
24 nameString + "\nSELECT DISTINCT ?person_uri ?comment {\n" +
25 " ?person_uri <http://xmlns.com/foaf/0.1/name> \"\" +
26 nameString + "\"@en .\n" +
27 " OPTIONAL { ?person_uri <http://www.w3.org/2000/01/rdf-schema#comment>\n" +
28 " ?comment . FILTER (lang(?comment) = 'en') } .\n" +
29 "}\nLIMIT 10\n"
30 }
31
32 public func get_SPARQL_for_PERSON_URI(aURI: String) -> String {
33 return
34 "# <" + aURI + ">\nSELECT DISTINCT ?comment (GROUP_CONCAT(DISTINCT ?birthpla\
35 ce; SEPARATOR=' | ') AS ?birthplace)\n (GROUP_CONCAT(DISTINCT ?almamater; SEPARATOR\
36 =' | ') AS ?almamater) (GROUP_CONCAT(DISTINCT ?spouse; SEPARATOR=' | ') AS ?spouse) \
37 {\n" +
38 " <" + aURI + "> <http://www.w3.org/2000/01/rdf-schema#comment> ?comment .\n" +
39 FILTER (lang(?comment) = 'en') .\n" +
40 " OPTIONAL { <" + aURI + "> <http://dbpedia.org/ontology/birthPlace> ?birth\
41 place } .\n" +
42 " OPTIONAL { <" + aURI + "> <http://dbpedia.org/ontology/almamater> ?almama\
43 ter } .\n" +

```

```

44 " OPTIONAL { <" + aURI + "> <http://dbpedia.org/ontology/spouse> ?spouse } \n
45 .\n" +
46 "} LIMIT 5\n"
47 }
48
49 public func get_display_text_for_PERSON_URI(personURI: String) -> [String] {
50 var ret: String = "\u2028(uri_to_display_text(uri: personURI))\n\n"
51 let person_details_sparql = get_SPARQL_for_PERSON_URI(aURI: personURI)
52 let person_details = sparqlDbPedia(query: person_details_sparql)
53
54 for pd in person_details {
55 //let comment = pd["comment"]
56 ret.append("\u2028(pd["comment"] ?? "")\n\n")
57 let subject_uris = pd["subject_uris"]
58 let uri_list: [String] = subject_uris?.components(separatedBy: " | ") ?? []
59 //ret.append("\n")
60 for u in uri_list {
61 let subject = uri_to_display_text(uri: u)
62 ret.append("\u2028(subject)\n")
63 //ret.append("\n")
64 if let spouse = pd["spouse"] {
65 if spouse.count > 0 {
66 ret.append("Spouse: \u2028(uri_to_display_text(uri: spouse))\n")
67 if let almamater = pd["almamater"] {
68 if almamater.count > 0 {
69 ret.append("Almamater: \u2028(uri_to_display_text(uri: almamater))\n")
70 if let birthplace = pd["birthplace"] {
71 if birthplace.count > 0 {
72 ret.append("Birthplace: \u2028(uri_to_display_text(uri: birthplace))\n")
73 }
74 }
75 return ["# SPARQL for a specific person:\n" + person_details_sparql, ret]
76 }
77
78 // "?place_uri <http://xmlns.com/foaf/0.1/name> \"\" + placeString + "\"@en .\n\n
79 " +
80
81 public func get_SPARQL_for_finding_URIs_for_PLACE_NAME(placeString: String)
82 -> String {
83 return
84 "# " + placeString + "\nSELECT DISTINCT ?place_uri ?comment {\n" +
85 " ?place_uri rdfs:label \"\" + placeString + "\"@en .\n" +
86 " ?place_uri <http://www.w3.org/1999/02/22-rdf-syntax-ns#type> <http://sche\
```

```

87 ma.org/Place> .\n" +
88 " OPTIONAL { ?place_uri <http://www.w3.org/2000/01/rdf-schema#comment>\n" +
89 " ?comment . FILTER (lang(?comment) = 'en') } .\n" +
90 "} LIMIT 10\n"
91 }
92
93 public func get_SPARQL_for_PLACE_URI(aURI: String) -> String {
94 return
95 "# <" + aURI + ">\nSELECT DISTINCT ?comment (GROUP_CONCAT(DISTINCT ?subject_\
96 uris; SEPARATOR=' | ') AS ?subject_uris) {\n" +
97 " <" + aURI + "> <http://www.w3.org/2000/01/rdf-schema#comment> ?comment .\n" +
98 FILTER (lang(?comment) = 'en') .\n" +
99 " OPTIONAL { <" + aURI + "> <http://purl.org/dc/terms/subject> ?subject_uri\\
100 s } .\n" +
101 "} LIMIT 5\n"
102 }
103
104 public func get_HTML_for_place_URI(placeURI: String) -> String {
105 var ret: String = "<h2>" + placeURI + "</h2>\n"
106 let place_details_sparql = get_SPARQL_for_PLACE_URI(aURI: placeURI)
107 let place_details = sparqlDbPedia(query: place_details_sparql)
108
109 for pd in place_details {
110 //let comment = pd["comment"]
111 ret.append("<p>\(pd["comment"] ?? "")</p>\n")
112 let subject_uris = pd["subject_uris"]
113 let uri_list: [String] = subject_uris?.components(separatedBy: " | ") ?? []
114 ret.append("\n")
115 for u in uri_list {
116 let subject = u.replacingOccurrences(of: "http://dbpedia.org/resource/Ca\
117 tegory:", with: "").replacingOccurrences(of: "_", with: " ").replacingOccurrences(of\
118 : "-", with: " ")
119 ret.append(" \(subject)\n")
120 }
121 ret.append("\n")
122 }
123 return ret
124 }
125
126 public func get_SPARQL_for_finding_URI_for_ORGANIZATION_NAME(orgString: String) -> \
127 String {
128 return
129 "# " + orgString + "\nSELECT DISTINCT ?org_uri ?comment {\n" +

```

```

130 " ?org_uri rdfs:label \"\" + orgString + "\"@en .\n" +
131 " ?org_uri <http://www.w3.org/1999/02/22-rdf-syntax-ns#type> <http://schema\
132 .org/Organization> .\n" +
133 " OPTIONAL { ?org_uri <http://www.w3.org/2000/01/rdf-schema#comment>\n" +
134 " ?comment . FILTER (lang(?comment) = 'en') } .\n" +
135 "} LIMIT 2\n"
136 }

```

The file **AppSparql** contains more utility functions for getting entity and relationship data from DBpedia:

```

1 // AppSparql.swift
2 // Created by ML Watson on 7/18/21.
3
4 import Foundation
5
6 let detailSparql = """
7 PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
8 select ?entity ?label ?description ?comment where {
9 ?entity rdfs:label "<name>"@en .
10 ?entity schema:description ?description . filter (lang(?description) = 'en') . f\
11 ilter(!regex(?description,"Wikimedia disambiguation page")) .
12 } limit 5000
13 """
14
15 let personSparql = """
16 select ?uri ?comment {
17 ?uri <http://xmlns.com/foaf/0.1/name> "<name>"@en .
18 ?uri <http://www.w3.org/2000/01/rdf-schema#comment> ?comment .
19 FILTER (lang(?comment) = 'en') .
20 }
21 """
22
23
24 let personDetailSparql = """
25 SELECT DISTINCT ?label ?comment
26
27 (GROUP_CONCAT (DISTINCT ?birthplace; SEPARATOR=' | ') AS ?birthplace)
28 (GROUP_CONCAT (DISTINCT ?almamater; SEPARATOR=' | ') AS ?almamater)
29 (GROUP_CONCAT (DISTINCT ?spouse; SEPARATOR=' | ') AS ?spouse) {
30 <name> <http://www.w3.org/2000/01/rdf-schema#comment> ?comment .
31 FILTER (lang(?comment) = 'en') .
32 OPTIONAL { <name> <http://dbpedia.org/ontology/birthPlace> ?birthplace } .

```

```

33 OPTIONAL { <name> <http://dbpedia.org/ontology/almaMater> ?almamater } .
34 OPTIONAL { <name> <http://dbpedia.org/ontology/spouse> ?spouse } .
35 OPTIONAL { <name> <http://www.w3.org/2000/01/rdf-schema#label> ?label } .
36 FILTER (lang(?label) = 'en') }
37 } LIMIT 10
38 """
39
40 let placeSparql = """
41 SELECT DISTINCT ?uri ?comment WHERE {
42 ?uri rdfs:label "<name>"@en .
43 ?uri <http://www.w3.org/2000/01/rdf-schema#comment> ?comment .
44 FILTER (lang(?comment) = 'en') .
45 ?place <http://www.w3.org/1999/02/22-rdf-syntax-ns#type> <http://schema.org/Place\
46 > .
47 } LIMIT 80
48 """
49
50 let organizationSparql = """
51 SELECT DISTINCT ?uri ?comment WHERE {
52 ?uri rdfs:label "<name>"@en .
53 ?uri <http://www.w3.org/2000/01/rdf-schema#comment> ?comment .
54 FILTER (lang(?comment) = 'en') .
55 ?uri <http://www.w3.org/1999/02/22-rdf-syntax-ns#type> <http://schema.org/Organiz\
56 ation> .
57 } LIMIT 80
58 """
59
60 func entityDetail(name: String) -> [Dictionary<String, String>] {
61 var ret: [Dictionary<String, String>] = []
62 let sparql = detailSparql.replacingOccurrences(of: "<name>", with: name)
63 print(sparql)
64 let r = sparqlDbPedia(query: sparql)
65 r.forEach { result in
66 print(result)
67 ret.append(result)
68 }
69 return ret
70 }
71
72 func personDetail(name: String) -> [Dictionary<String, String>] {
73 var ret: [Dictionary<String, String>] = []
74 let sparql = personSparql.replacingOccurrences(of: "<name>", with: name)
75 print(sparql)

```

```
76 let r = sparqlDbPedia(query: sparql)
77 r.forEach { result in
78 print(result)
79 ret.append(result)
80 }
81 return ret
82 }
83
84 func placeDetail(name: String) -> [Dictionary<String, String>] {
85 var ret: [Dictionary<String, String>] = []
86 let sparql = placeSparql.replacingOccurrences(of: "<name>", with: name)
87 print(sparql)
88 let r = sparqlDbPedia(query: sparql)
89 r.forEach { result in
90 print(result)
91 ret.append(result)
92 }
93 return ret
94 }
95
96 func organizationDetail(name: String) -> [Dictionary<String, String>] {
97 var ret: [Dictionary<String, String>] = []
98 let sparql = organizationSparql.replacingOccurrences(of: "<name>", with: name)
99 print(sparql)
100 let r = sparqlDbPedia(query: sparql)
101 r.forEach { result in
102 print(result)
103 ret.append(result)
104 }
105 return ret
106 }
107
108 public func processEntities(inputString: String) -> [(name: String, type: String, ur\
109 i: String, comment: String)] {
110 let entities = getEntities(text: inputString)
111 var augmentedEntities: [(name: String, type: String, uri: String, comment: Strin\
112 g)] = []
113 for (entityName, entityType) in entities {
114 print("** entityName:", entityName, "entityType:", entityType)
115 if entityType == "PersonalName" {
116 let data = personDetail(name: entityName)
117 for d in data {
118 augmentedEntities.append((name: entityName, type: entityType,
```

```

119 uri: "<" + d["uri"]! + ">", comment: "<" + d["comment"]! + ">"))
120 }
121 }
122 if entityType == "OrganizationName" {
123 let data = organizationDetail(name: entityName)
124 for d in data {
125 augmentedEntities.append((name: entityName, type: entityType,
126 uri: "<" + d["uri"]! + ">", comment: "<" + d["comment"]! + ">"))
127 }
128 }
129 if entityType == "PlaceName" {
130 let data = placeDetail(name: entityName)
131 for d in data {
132 augmentedEntities.append((name: entityName, type: entityType,
133 uri: "<" + d["uri"]! + ">", comment: "<" + d["comment"]! + ">"))
134 }
135 }
136 }
137 return augmentedEntities
138 }
139
140
141 extension Array where Element: Hashable {
142 func uniqueValuesHelper() -> [Element] {
143 var addedDict = [Element: Bool]()
144 return filter { addedDict.updateValue(true, forKey: $0) == nil }
145 }
146 mutating func uniqueValues() {
147 self = self.uniqueValuesHelper()
148 }
149 }
150
151
152 func getAllRelationships(inputString: String) -> [String] {
153 let augmentedEntities = processEntities(inputString: inputString)
154 var relationshipTriples: [String] = []
155 for ae1 in augmentedEntities {
156 for ae2 in augmentedEntities {
157 if ae1 != ae2 {
158 let er1 = dbpediaGetRelationships(entity1Uri: ae1.uri,
159 entity2Uri: ae2.uri)
160 relationshipTriples.append(contentsOf: er1)
161 let er2 = dbpediaGetRelationships(entity1Uri: ae2.uri,

```

```

162 entity2Uri: ae1.uri)
163 relationshipTriples.append(contentsOf: er2)
164 }
165 }
166 }
167 relationshipTriples.uniqueValues()
168 relationshipTriples.sort()
169 return relationshipTriples
170 }
```

## AppleBERT

The files in the directory AppleBERT were copied from Apple's example [https://developer.apple.com/documentation/integration\\_samples/finding\\_answers\\_to\\_questions\\_in\\_a\\_text\\_document<sup>73</sup>](https://developer.apple.com/documentation/integration_samples/finding_answers_to_questions_in_a_text_document<sup>73</sup>) with a few changes to get returned results in a convenient format for this application. Apple's BERT documentation is excellent and you should review it.

## Relationships

The file Relationships.swift fetches relationship data for pairs of DBpedia entities. Note that the first SPARQL template has variable slots <e1> and <e2> that are replaced at runtime with URIs representing the entities that we are searching for relationships between these two entities:

```

1 // relationships between DBpedia entities
2
3 let relSparql = """
4 SELECT DISTINCT ?p {<e1> ?p <e2> .FILTER (!regex(str(?p), 'wikiPage', 'i'))} LIMIT 5
5 """
6
7 public func dbpediaGetRelationships(entity1Uri: String, entity2Uri: String)
8 -> [String] {
9 var ret: [String] = []
10 let sparql1 = relSparql.replacingOccurrences(of: "<e1>",
11 with: entity1Uri).replacingOccurrences(of: "<e2>",
12 with: entity2Uri)
13 let r1 = sparqlDbPedia(query: sparql1)
14 r1.forEach { result in
15 if let relName = result["p"] {
16 let rdfStatement = entity1Uri + " <" + relName + "> " + entity2Uri + "."
17 print(rdfStatement)
18 }
19 }
20 }
```

---

<sup>73</sup>[https://developer.apple.com/documentation/coreml/model\\_integration\\_samples/finding\\_answers\\_to\\_questions\\_in\\_a\\_text\\_document](https://developer.apple.com/documentation/coreml/model_integration_samples/finding_answers_to_questions_in_a_text_document)

```
18 ret.append(rdfStatement)
19 }
20 }
21 let sparql2 = relSparql.replacingOccurrences(of: "<e1>",
22 with: entity2Uri).replacingOccurrences(of: "<e2>",
23 with: entity1Uri)
24 let r2 = sparqlDbPedia(query: sparql2)
25 r2.forEach { result in
26 if let relName = result["p"] {
27 let rdfStatement = entity2Uri + " <" + relName + "> " + entity1Uri + "."
28 print(rdfStatement)
29 ret.append(rdfStatement)
30 }
31 }
32 return Array(Set(ret))
33 }
34
35 public func uriToPrintName(_ uri: String) -> String {
36 let slashIndex = uri.lastIndex(of: "/")
37 if slashIndex == nil { return uri }
38 var s = uri[slashIndex!...]
39 s = s.dropFirst()
40 if s.count > 0 { s.removeLast() }
41 return String(s).replacingOccurrences(of: "_", with: " ")
42 }
43
44 public func relationshipsInEnglish(rs: [String]) -> String {
45 var lines: [String] = []
46 for r in rs {
47 let triples = r.split(separator: " ", maxSplits: 3,
48 omittingEmptySubsequences: true)
49 if triples.count > 2 {
50 lines.append(uriToPrintName(String(triples[0])) + " " +
51 uriToPrintName(String(triples[1])) + " " +
52 uriToPrintName(String(triples[2])))
53 } else {
54 lines.append(r)
55 }
56 }
57 let linesNoDuplicates = Set(lines)
58 return linesNoDuplicates.joined(separator: "\n")
59 }
```

## NLP

The file **NlpWhiteboard** provides high level NLP utility functions for the application:

```

1 //
2 // NlpWhiteboard.swift
3 // KGN
4 //
5 // Copyright © 2021 Mark Watson. All rights reserved.
6 //
7
8 public struct NlpWhiteboard {
9
10 var originalText: String = ""
11 var people: [String] = []
12 var places: [String] = []
13 var organizations: [String] = []
14 var sparql: String = ""
15
16 init() { }
17
18 mutating func set_text(originalText: String) {
19 self.originalText = originalText
20 let (people, places, organizations) = getAllEntities(text: originalText)
21 self.people = people; self.places = places; self.organizations = organizatio\
22 ns
23 }
24
25 mutating func query_to_choices(behindTheScenesSparqlText: inout String)
26 -> [[[String]]] { // return inner: [comment, uri]
27 var ret: Set<[String]> = []
28 if people.count > 0 {
29 for i in 0...(people.count - 1) {
30 self.sparql =
31 get_SPARQL_for_finding_URI_for_PERSON_NAME(nameString: people[i])
32 behindTheScenesSparqlText += self.sparql
33 let results = sparqlDbPedia(query: self.sparql)
34 if results.count > 0 {
35 ret.insert(results.map { [($0["comment"]
36 ?? ""),
37 ($0["person_uri"] ?? "")] })
38 }
39 }
40 }
41 }
42 }
```

```

40 }
41 if organizations.count > 0 {
42 for i in 0...(organizations.count - 1) {
43 self.sparql = get_SPARQL_for_finding_URIIs_for_ORGANIZATION_NAME(
44 orgString: organizations[i])
45 behindTheScenesSparqlText += self.sparql
46 let results = sparqlDbPedia(query: self.sparql)
47 if results.count > 0 {
48 ret.insert(results.map { [($0["comment"] ?? "")], ($0["org_uri"] ?? "")] })
49 }
50 }
51 }
52 if places.count > 0 {
53 for i in 0...(places.count - 1) {
54 self.sparql = get_SPARQL_for_finding_URIIs_for_PLACE_NAME(
55 placeString: places[i])
56 behindTheScenesSparqlText += self.sparql
57 let results = sparqlDbPedia(query: self.sparql)
58 if results.count > 0 {
59 ret.insert(results.map { [($0["comment"] ?? "")], ($0["place_uri"] ?? "")] })
60 }
61 }
62 }
63 //print("\n\n++++++ ret:\n", ret, "\n\n")
64 return Array(ret)
65 }
66 }
67 }
68 }
```

The file **NLUtils.swift** provides lower level NLP utilities:

```

1 // NLUtils.swift
2 // KGN
3 //
4 // Copyright © 2021 Mark Watson. All rights reserved.
5 //
6
7 import Foundation
8 import NaturalLanguage
9
10 public func getPersonDescription(personName: String) -> [String] {
11 let sparql = get_SPARQL_for_finding_URIIs_for_PERSON_NAME(nameString: personName)
```

```
12 let results = sparqlDbPedia(query: sparql)
13 return [sparql, results.map {
14 ($0["comment"] ?? $0["abstract"] ?? "") }.joined(separator: " . ")]
15 }
16
17
18 public func getPlaceDescription(placeName: String) -> [String] {
19 let sparql = get_SPARQL_for_finding_URIs_for_PLACE_NAME(placeString: placeName)
20 let results = sparqlDbPedia(query: sparql)
21 return [sparql, results.map { ($0["comment"] ???
22 $0["abstract"] ?? "") }.joined(separator: " . ")]
23 }
24
25 public func getOrganizationDescription(organizationName: String) -> [String] {
26 let sparql = get_SPARQL_for_finding_URIs_for_ORGANIZATION_NAME(
27 orgString: organizationName)
28 let results = sparqlDbPedia(query: sparql)
29 print("== getOrganizationDescription results =\n", results)
30 return [sparql, results.map { ($0["comment"] ?? $0["abstract"] ?? "") }
31 .joined(separator: " . ")]
32 }
33
34 let tokenizer = NLTokenizer(unit: .word)
35 let tagger = NSLinguisticTagger(tagSchemes:[.tokenType, .language, .lexicalClass,
36 .nameType, .lemma], options: 0)
37 let options: NSLinguisticTagger.Options =
38 [.omitPunctuation, .omitWhitespace, .joinNames]
39
40 let tokenizerOptions: NSLinguisticTagger.Options =
41 [.omitPunctuation, .omitWhitespace, .joinNames]
42
43 public func getEntities(text: String) -> [(String, String)] {
44 var words: [(String, String)] = []
45 tagger.string = text
46 let range = NSRange(location: 0, length: text.utf16.count)
47 tagger.enumerateTags(in: range, unit: .word,
48 scheme: .nameType, options: options) { tag, tokenRange, stop in
49 let word = (text as NSString).substring(with: tokenRange)
50 let tagType = tag?.rawValue ?? "unkown"
51 if tagType != "unkown" && tagType != "OtherWord" {
52 words.append((word, tagType))
53 }
54 }
```

```
55 return words
56 }
57
58 public func tokenizeText(text: String) -> [String] {
59 var tokens: [String] = []
60 tokenizer.string = text
61 tokenizer.enumerateTokens(in: text.startIndex..
```

```
98 func splitLongStrings(_ s: String, limit: Int) -> String {
99 var ret: [String] = []
100 let tokens = s.split(separator: " ")
101 var subLine = ""
102 for token in tokens {
103 if subLine.count > limit {
104 ret.append(subLine)
105 subLine = ""
106 } else {
107 subLine = subLine + " " + token
108 }
109 }
110 if subLine.count > 0 {
111 ret.append(subLine)
112 }
113 return ret.joined(separator: "\n")
114 }
```

## Views

This is not a book about SwiftUI programming, and indeed I expect many of you dear readers know much more about UI development with SwiftUI than I do. I am not going to list the four view files:

- MainView.swift
- QueryView.swift
- AboutView.swift
- InfoView.swift

## Main KGN

The top level app code in the file **KGNApp.swift** is fairly simple. I hardcoded the window size for macOS and the window sizes for running this example on iPadOS or iOS are commented out:

```
1 import SwiftUI
2
3 @main
4 struct KGNApp: App {
5 var body: some Scene {
6 WindowGroup {
7 MainView()
8 .frame(width: 1200, height: 770) // << here !!
9 // .frame(width: 660, height: 770) // << here !!
10 // .frame(width: 500, height: 800) // << here !!
11 }
12 }
13 }
```

I was impressed by the SwiftUI framework. Applications are fairly portable across macOS, iOS, and iPadOS. I am not a UI developer by profession (as this application shows) but I enjoyed learning just enough about SwiftUI to write this example application.

# **Part 5: Apple Intelligence**

## **Developers Can Now Weave Apple Intelligence Directly Into Their Apps**

Dear reader, Apple has opened up its new “Apple Intelligence” system to third-party developers, providing a suite of tools and APIs to create more personalized and powerful applications. Developers can primarily leverage this technology through the new Foundation Models framework, which offers direct access to Apple’s on-device generative models. This allows for the integration of sophisticated features like text summarization, content creation, and message rewriting directly within an app’s interface, all while processing user data locally on the device to ensure privacy.

Furthermore, developers can utilize an expanded App Intents framework to make their application’s content and functionalities accessible to system-wide services like Siri and Shortcuts. This deeper integration allows users to interact with apps using natural language and create complex, multi-app workflows. For instance, a user could ask Siri to “find all photos of my dog from last summer and create a collage in my favorite editing app,” and the system, through App Intents, would understand and execute this command. The Image Playground API also allows for the seamless inclusion of image generation capabilities, enabling users to create unique visuals within the context of the app they are using.

### **Key Advantages for Developers:**

- Enhanced User Experience: By integrating features like personalized content suggestions, intelligent summarization, and natural language interaction, developers can create more intuitive and engaging applications that anticipate user needs.
- On-Device Processing for Privacy and Performance: A significant advantage is the on-device nature of Apple Intelligence. This approach enhances user privacy and security by keeping personal data on the device. It also improves app performance and responsiveness by reducing reliance on cloud-based processing.
- Access to Powerful Generative Models: Developers gain access to Apple’s sophisticated generative models without the need to build and train their own. This allows for the rapid implementation of advanced AI features that were previously complex and resource-intensive to develop.
- Deeper System Integration: The enhanced App Intents framework provides a new level of integration with core iOS, iPadOS, and macOS features. This allows apps to become more discoverable and useful within the broader Apple ecosystem, extending their functionality beyond the confines of the app itself.

# Using Apple Intelligence's Default System Model To Build a Chat Command Line Tool

Here we create a simple command line chat tool. The LLM specific code can also be used in iOS, iPadOS, and macOS applications. This chat tool will use the local system LLM for simple queries and will transparently call a more capable model on Apple's servers in a secure and privacy preserving sandbox.

The following Swift package manifest defines an executable command-line program named **chattool** designed to run on macOS. The tool relies on Apple's swift-argument-parser package to process command-line arguments.

Package.swift:

```
1 // swift-tools-version: 6.2
2 import PackageDescription
3
4 let package = Package(
5 name: "ChatTool",
6 platforms: [.macOS(.v26)], // macOS 15/16 is fine too
7 products: [
8 .executable(name: "chattool", targets: ["ChatTool"])
9],
10 dependencies: [
11 // □ Apple's official argument-parsing package
12 .package(url: "https://github.com/apple/swift-argument-parser.git",
13 from: "1.3.0")
14],
15 targets: [
16 .executableTarget(
17 name: "ChatTool",
18 // expose the ArgumentParser product to the target
19 dependencies: [
20 .product(name: "ArgumentParser", package: "swift-argument-parser")
21 // FoundationModels is an Apple framework, no SPM entry needed
22],
23 path: "Sources/ChatTool" // adjust if your path differs
```

```
24)
25]
26 }
```

### ChatTool.swift:

This Swift code implements a command-line chat interface that interacts directly with Apple's native FoundationModels framework. Upon launch, it verifies the system's default language model is available, initializes a LanguageModelSession with a hard-coded system prompt and temperature, and then enters a read-evaluate-print loop. The program continuously accepts user input from the console and sends it to the language model for processing until the user quits.

The core of the implementation leverages modern Swift concurrency (async/await) to handle the model's output as an asynchronous stream. For each user prompt, it iterates through the text fragments as they are generated by session.streamResponse, writing only the new characters to standard output to create a real-time, typewriter-like effect. To enhance usability, it uses a DispatchSource to set up a signal handler for SIGINT (Ctrl+C), allowing a user to cancel the current in-progress stream from the model without terminating the entire chat application.

```
1 import Foundation
2 import FoundationModels
3 import Dispatch
4
5 @main
6 struct ChatCLI {
7 static func main() async throws {
8 // Hard-coded defaults
9 let temperature = 0.2
10 let sysPrompt = "You are a helpful assistant."
11
12 // Verify model
13 let model = SystemLanguageModel.default
14 guard model.isAvailable else {
15 throw RuntimeError("Model unavailable: \(model.availability)")
16 }
17
18 let session = LanguageModelSession(instructions: sysPrompt)
19 print("Temperature: \(temperature)")
20 print("System Prompt: \(sysPrompt)")
21 let options = GenerationOptions(temperature: temperature)
22
23 print("Apple-Intelligence chat (streaming, T=0.2). Type /quit to exit.\n")
24
25 while true {
```

```

26 print("Enter your message: ", terminator: "")
27 guard let prompt = readLine(strippingNewline: true) else { break }
28 if prompt.isEmpty || prompt == "/quit" { break }
29
30 var previous = "" // text already printed
31
32 let task = Task {
33 for try await part in session.streamResponse(to: prompt, options: op\
34 tions) {
35 let delta = part.dropFirst(previous.count) // new characters only
36 if !delta.isEmpty {
37 FileHandle.standardOutput.write(Data(delta.utf8))
38 fflush(stdout)
39 previous = part
40 }
41 }
42 print() // newline when complete
43 }
44
45 // ^C cancels the streaming task
46 signal(SIGINT, SIG_IGN)
47 let sigSrc = DispatchSource.makeSignalSource(signal: SIGINT, queue: .mai\
48 n)
49 sigSrc.setEventHandler { task.cancel() }
50 sigSrc.resume()
51 defer { sigSrc.cancel() }
52
53 _ = try await task.value
54 }
55 }
56 }
57
58 /// Simple error wrapper
59 struct RuntimeError: Error, CustomStringConvertible {
60 let description: String
61 init(_ msg: String) { description = msg }
62 }
```

Here are the first few lines of output given the prompt *Describe the math for calculating the orbit of Jupiter, then write a very short design for a Python script:*

```
1 $ swift run
2 [1/1] Planning build
3 Building for debugging...
4 [1/1] Write swift-version-39B54973F684ADAB.txt
5 Build of product 'chattool' complete! (0.16s)
6 Temperature: 0.2
7 System Prompt: You are a helpful assistant.
8 Apple-Intelligence chat (streaming, T=0.2). Type /quit to exit.
9
10 Enter your message: Describe the math for calculating the orbit of Jupiter, then wri\
11 te a very short design for a Python script
12 Calculating the orbit of Jupiter involves solving Kepler's laws of planetary motion, \
13 which describe the elliptical orbits of planets around the Sun. The key equations i\
14 nvolve gravitational forces and conservation laws. Here's a brief overview of the ma\
15 th involved:
16
17 ### Key Concepts:
18
19 1. **Kepler's Laws:**
20 - **First Law (Law of Ellipses):** Planets move along ellipses with the Sun at on\
21 e focus.
22 - **Second Law (Law of Equal Areas):** A line segment joining a planet and the Su\
23 n sweeps out equal areas during equal intervals of time.
24 - **Third Law (Law of Harmonies):** The square of the orbital period ($\langle T \rangle$) is p\
25 roportional to the cube of the semi-major axis ($\langle a \rangle$): $\langle T^2 = \frac{4\pi^2}{GM}a^3 \rangle$.
26
27
28 2. **Gravitational Force:**
29 - The gravitational force between two masses ($\langle m_1 \rangle$ and $\langle m_2 \rangle$) is given by Ne\
30 wton's law: $\langle F = G \frac{m_1 m_2}{r^2} \rangle$, where $\langle G \rangle$ is the gravitational constant\
31 and $\langle r \rangle$ is the distance between centers.
32
33 3. **Centripetal Force:**
34 - For circular orbits, centripetal force equals gravitational force: $\langle F = \frac{mv^2}{r} \rangle$.
```

# Using Apple Intelligence's Default System Model To Build a Coding Assistant Command Line Tool

This tool looks in the current directory and all subdirectories for source code files and describes them and then enters a chat loop for talking about the code.

**Package.swift:**

```
1 // swift-tools-version: 6.2
2 // The swift-tools-version declares the minimum version of Swift required to build t\
3 his package.
4
5 import PackageDescription
6
7 let package = Package(
8 name: "CodingCLI",
9
10 // 1 Tell SwiftPM we require at least macOS 12 so
11 // `Task.value`, async/await, and FoundationModels are available.
12 platforms: [
13 .macOS(.v26)
14],
15
16 products: [
17 .executable(name: "CodingCLI", targets: ["CodingCLI"])
18],
19
20 targets: [
21 .executableTarget(
22 name: "CodingCLI",
23
24 // 2 Link the system framework that ships with Xcode 17+
25 // (no external dependency required).
26 linkerSettings: [
27 .linkedFramework("FoundationModels")
28]
29)

```

```
30]
31)
```

### CodingCLI.swift:

```
1 import Foundation
2 import FoundationModels
3 import Dispatch
4
5 @main
6 struct CodingCLI {
7 static func main() async throws {
8 // ---- 1. Gather candidate source files ----
9 let exts = ["swift", "py", "lisp"]
10 var blobs: [String] = []
11
12 let enumerator = FileManager.default.enumerator(atPath: ".")!
13
14 while let path = enumerator.nextObject() as? String { // avoids @noasync
15 guard let ext = path.split(separator: ".").last,
16 exts.contains(ext.lowercased()) else { continue }
17
18 if let data = FileManager.default.contents(atPath: path),
19 data.count < 8 * 1024 { // keep size filter
20 let text = String(decoding: data, as: UTF8.self) // non-optional
21 blobs.append("### \(path) ###\n\(text)")
22 }
23 }
24
25 let doc = blobs.joined(separator: "\n")
26 let summary = try await Self.summarize(doc)
27 print("\n==== Project Summary ====\n\(summary)\n")
28
29 // ---- 2. Start interactive chat loop ----
30 let session = LanguageModelSession(instructions:
31 "You are a helpful assistant.")
32 let options = GenerationOptions(temperature: 0.2)
33 print("Apple-Intelligence chat (streaming, T=0.2). Type /quit to exit.\n")
34
35 while let prompt = readLine(strippingNewline: true) {
36 if prompt.isEmpty || prompt == "/quit" { break }
37
38 var printed = "
```

```

39 let task = Task {
40 for try await part in session.streamResponse(to: prompt,
41 options: options) {
42 let delta = part.dropFirst(printed.count)
43 if !delta.isEmpty {
44 FileHandle.standardOutput.write(Data(delta.utf8))
45 fflush(stdout)
46 printed = part
47 }
48 }
49 print()
50 }
51
52 signal(SIGINT, SIG_IGN)
53 let sig = DispatchSource.makeSignalSource(signal: SIGINT, queue: .main)
54 sig.setEventHandler { task.cancel() }
55 sig.resume()
56 defer { sig.cancel() }
57
58 _ = try await task.value
59 }
60 }
61
62 // ---- 3. Helper: summarize all code ----
63 static func summarize(_ text: String) async throws -> String {
64 let session = LanguageModelSession(
65 instructions: """
66 Summarize the following multi-file project. \
67 For each file give one bullet explaining its role, then \
68 a two-sentence overall description.
69 """
70)
71 let prompt = text.prefix(24 * 1024) // safety window
72 let resp = try await
73 session.respond(to: String(prompt),
74 options: GenerationOptions(temperature: 0))
75 return resp.content // unwrap Response<String>
76 }
77 }
```

Here is the output for running this tool in its own source directory:

```
1 $ swift run
2 Building for debugging...
3 [8/8] Applying CodingCLI
4 Build of product 'CodingCLI' complete! (3.23s)
5
6 === Project Summary ===
7 ### test.py
8 - **Role:** This script interacts with Groq to perform a chat completion task.
9 - **Description:** It sets up a chat session using Groq's API, sends a specific message, and prints the response, showcasing how to utilize Groq for conversational AI tasks.
10
11
12 ### Package.swift
13 - **Role:** Defines the Swift package configuration for the CodingCLI project.
14 - **Description:** This file specifies the project's platform requirements, defines the executable product, and outlines the executable target with necessary dependencies.
15
16
17
18 ### Sources/CodingCLI/CodingCLI.swift
19 - **Role:** Serves as the entry point for the CodingCLI application, handling file summarization and chat interaction.
20 - **Description:** It processes source files to generate a summary, and manages an interactive chat loop using a language model, demonstrating integration of summarization and conversational AI within a Swift package.
21
22
23
24
25
26 Apple-Intelligence chat (streaming, T=0.2). Type /quit to exit.
```

# **Book Wrap Up**

I hope that you dear reader enjoyed this short book. While I enjoy programming in Swift and appreciate how well Apple has integrated machine learning capabilities in their iOS/iPadOS/macOS ecosystems, I still find myself writing most of my experimental code in Lisp languages and using Python for deep learning experiments and projects. That said, I am very happy that I have done the work to add Swift, CoreML, and SwiftUI to my personal programming tool belt.

I usually update my eBooks so if there is some topic or application domain that you would like added to future versions of this book, then please let me know. My email address is markw <at> markwatson <dot> com.