# Deep dive into coroutines
## on the example of Kotlin implementation

Maciej Procyk



March 27, 2023

Introduction
000000000000

Coroutines in details
00000000000000000000

Higher-level APIs
00000000

Appendix
00

# Table of contents

**K** Kotlin

# Introduction

**Introduction**
○●○○○○○○○○○○○

Coroutines in details
○○○○○○○○○○○○○○○○○○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

Main ideas

# General description

- described as "functions whose execution you can pause"

**K** Kotlin

# General description

- described as "functions whose execution you can pause"
- computation can be paused and resumed without blocking a thread

**K** Kotlin

# General description

- described as "functions whose execution you can pause"
- computation can be paused and resumed without blocking a thread
- generalization of subroutines for cooperative multitasking

**K** Kotlin

# General description

- described as "functions whose execution you can pause"
- computation can be paused and resumed without blocking a thread
- generalization of subroutines for cooperative multitasking
- fast context switching

**K** Kotlin

# General description

- described as "functions whose execution you can pause"
- computation can be paused and resumed without blocking a thread
- generalization of subroutines for cooperative multitasking
- fast context switching
- perfect tool for implementing iterators, infinite lists, pipes

**K** Kotlin

# General description

- described as "functions whose execution you can pause"
- computation can be paused and resumed without blocking a thread
- generalization of subroutines for cooperative multitasking
- fast context switching
- perfect tool for implementing iterators, infinite lists, pipes
- asynchronous code in synchronous manner

**K** Kotlin

# Fundamental characteristics

1 the values of data local to a coroutine persist between
  successive calls

**K** Kotlin

# Fundamental characteristics

1. the values of data local to a coroutine persist between successive calls

2. the execution of a coroutine is suspended as control leaves it, only to carry on where it left off when control re-enters the coroutine at some later stage

**K** Kotlin

# Asynchronous programming

- nowadays, we write the programs that wait for a lot of time

**K** Kotlin

Introduction
○○○●○○○○○○○○○
Coroutines in details
○○○○○○○○○○○○○○○○○○○○○○
Higher-level APIs
○○○○○○○○
Appendix
○○

Practical background

# Asynchronous programming

- nowadays, we write the programs that wait for a lot of time
- we write them using threads, but there are a few limitations

**K** Kotlin

Introduction
○○○●○○○○○○○○○
Coroutines in details
○○○○○○○○○○○○○○○○○○○○○○
Higher-level APIs
○○○○○○○○
Appendix
○○

Practical background

# Asynchronous programming

- nowadays, we write the programs that wait for a lot of time
- we write them using threads, but there are a few limitations
- the number of threads cannot be too high, as they consume a lot of memory

**K** Kotlin

Introduction
○○○●○○○○○○○○○

Coroutines in details
○○○○○○○○○○○○○○○○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

Practical background

# Asynchronous programming

- nowadays, we write the programs that wait for a lot of time
- we write them using threads, but there are a few limitations
- the number of threads cannot be too high, as they consume a lot of memory
- e.g. on Linux x64 default stack size (configured with -XX:ThreadStackSize or -Xss option) is 256 KB

**K** Kotlin

Introduction
0000●000000000
Practical background

Coroutines in details
0000000000000000000000

Higher-level APIs
00000000

Appendix
00

# Asynchronous programming

- nowadays, we write the programs that wait for a lot of time
- we write them using threads, but there are a few limitations
- the number of threads cannot be too high, as they consume a lot of memory
- e.g. on Linux x64 default stack size (configured with `-XX:ThreadStackSize` or `-Xss` option) is 256 KB
- which means that 100K threads would consume about 24 GB of memory

**K** Kotlin

# Asynchronous programming

- nowadays, we write the programs that wait for a lot of time
- we write them using threads, but there are a few limitations
- the number of threads cannot be too high, as they consume a lot of memory
- e.g. on Linux x64 default stack size (configured with -XX:ThreadStackSize or -Xss option) is 256 KB
- which means that 100K threads would consume about 24 GB of memory
- that's a lot and what if we think about 1M threads

**K** Kotlin

Introduction
○○○○●○○○○○○○○○
Coroutines in details
○○○○○○○○○○○○○○○○○○○○○○
Higher-level APIs
○○○○○○○○
Appendix
○○

Practical background

# Coroutines vs threads

- the difference is that threads are scheduled on cores preemptively, while coroutines are scheduled onto threads cooperatively

**K** Kotlin

Introduction
○○○○●○○○○○○○○
Coroutines in details
○○○○○○○○○○○○○○○○○○○○○○
Higher-level APIs
○○○○○○○○
Appendix
○○
Practical background

# Coroutines vs threads

- the difference is that threads are scheduled on cores preemptively, while coroutines are scheduled onto threads cooperatively

- starting a thread on JVM requires starting native thread, which consumes memory for its stack

**K** Kotlin

Introduction
○○○○○●○○○○○○○○

Coroutines in details
○○○○○○○○○○○○○○○○○○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

Practical background

## Coroutines vs threads

- the difference is that threads are scheduled on cores preemptively, while coroutines are scheduled onto threads cooperatively
- starting a thread on JVM requires starting native thread, which consumes memory for its stack
- switching between threads requires going through system's kernel

**K** Kotlin

Introduction
○○○○●○○○○○○○○○
Practical background

Coroutines in details
○○○○○○○○○○○○○○○○○○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

# Coroutines vs threads

- the difference is that threads are scheduled on cores preemptively, while coroutines are scheduled onto threads cooperatively
- starting a thread on JVM requires starting native thread, which consumes memory for its stack
- switching between threads requires going through system's kernel
- and is expensive in terms of CPU cycles consumed during that operation

**K** Kotlin

Introduction
○○○○○○●○○○○○○
Coroutines in details
○○○○○○○○○○○○○○○○○○○○○○○
Higher-level APIs
○○○○○○○○
Appendix
○○
Practical background

# Coroutines vs threads

- coroutine is user-level abstraction (implemented by the language compiler in case of Kotlin)

**K** Kotlin

Introduction
○○○○○○●○○○○○○○
Practical background

Coroutines in details
○○○○○○○○○○○○○○○○○○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

# Coroutines vs threads

- coroutine is user-level abstraction (implemented by the language compiler in case of Kotlin)
- in the simplest case it's a single reference object in JVM heap memory

**K** Kotlin

# Coroutines vs threads

- coroutine is user-level abstraction (implemented by the language compiler in case of Kotlin)
- in the simplest case it's a single reference object in JVM heap memory
- switching between coroutines is as simple (and cheap) as invoking a regular function

**K** Kotlin

Introduction
○○○○○○○●○○○○○○○
Coroutines in details
○○○○○○○○○○○○○○○○○○○○○○○○
Higher-level APIs
○○○○○○○○
Appendix
○○

Practical background

# Sample code starting 100 000 threads

```kotlin
 1    fun main() {
 2        val counter = AtomicInteger(0)
 3        List(100_000) {
 4            thread {
 5                sleep(1_000)
 6                counter.incrementAndGet()
 7            }
 8        }.forEach { it.join() }
 9        println("counter = ${counter.get()}")
10    }
```

**K** Kotlin

Introduction
○○○○○○○●○○○○○○

Coroutines in details
○○○○○○○○○○○○○○○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

Practical background

# Sample code starting 100 000 threads

```
 1   fun main() {
 2       val counter = AtomicInteger(0)
 3       List(100_000) {
 4           thread {
 5               sleep(1_000)
 6               counter.incrementAndGet()
 7           }
 8       }.forEach { it.join() }
 9       println("counter = ${counter.get()}")
10   }
```

- possible, but takes a lot of time (and actually they don't exist in the same time)

**K** Kotlin

Introduction
○○○○○○○●○○○○○○
Practical background

Coroutines in details
○○○○○○○○○○○○○○○○○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

# Sample code starting 100 000 threads

```
 1   fun main() {
 2       val counter = AtomicInteger(0)
 3       List(100_000) {
 4           thread {
 5               sleep(1_000)
 6               counter.incrementAndGet()
 7           }
 8       }.forEach { it.join() }
 9       println("counter = ${counter.get()}")
10   }
```

- possible, but takes a lot of time (and actually they don't exist in the same time)

- if we increase the sleep time the number of threads existing in the same time would increase

**K** Kotlin

Introduction
○○○○○○○●○○○○○○

Coroutines in details
○○○○○○○○○○○○○○○○○○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

Practical background

# Sample code starting 100 000 threads

```
1   fun main() {
2       val counter = AtomicInteger(0)
3       List(100_000) {
4           thread {
5               sleep(1_000)
6               counter.incrementAndGet()
7           }
8       }.forEach { it.join() }
9       println("counter = ${counter.get()}")
10  }
```

- possible, but takes a lot of time (and actually they don't exist in the same time)
- if we increase the sleep time the number of threads existing in the same time would increase
- and we would end up with Exception in thread "main" java.lang.OutOfMemoryError ...

**K** Kotlin

Introduction
○○○○○○○●○○○○○
Practical background

Coroutines in details
○○○○○○○○○○○○○○○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

# Sample code starting 100 000 jobs

```
 1    fun main() = runBlocking {
 2        val counter = AtomicInteger(0)
 3        List(100_000) {
 4            launch {
 5                delay(1_000)
 6                counter.incrementAndGet()
 7            }
 8        }.forEach { it.join() }
 9        println("counter = ${counter.get()}")
10    }
```

**K** Kotlin

Introduction
○○○○○○○●○○○○○

Coroutines in details
○○○○○○○○○○○○○○○○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

Practical background

# Sample code starting 100 000 jobs

```
 1   fun main() = runBlocking {
 2       val counter = AtomicInteger(0)
 3       List(100_000) {
 4           launch {
 5               delay(1_000)
 6               counter.incrementAndGet()
 7           }
 8       }.forEach { it.join() }
 9       println("counter = ${counter.get()}")
10   }
```

- explicitly defined context is needed to start the whole process

**K** Kotlin

Introduction
000000000000000
Coroutines in details
0000000000000000000000
Higher-level APIs
00000000
Appendix
00

Practical background

# Sample code starting 100 000 jobs

```
 1   fun main() = runBlocking {
 2       val counter = AtomicInteger(0)
 3       List(100_000) {
 4           launch {
 5               delay(1_000)
 6               counter.incrementAndGet()
 7           }
 8       }.forEach { it.join() }
 9       println("counter = ${counter.get()}")
10   }
```

- explicitly defined context is needed to start the whole process
- code runs over 10 times faster than the threads implementation

**K** Kotlin

Introduction
0000000●00000

Coroutines in details
0000000000000000000000

Higher-level APIs
00000000

Appendix
00

Practical background

# Sample code starting 100 000 jobs

```kotlin
 1   fun main() = runBlocking {
 2       val counter = AtomicInteger(0)
 3       List(100_000) {
 4           launch {
 5               delay(1_000)
 6               counter.incrementAndGet()
 7           }
 8       }.forEach { it.join() }
 9       println("counter = ${counter.get()}")
10   }
```

- explicitly defined context is needed to start the whole process
- code runs over 10 times faster than the threads implementation
- because it uses computer resources more thoughtful

**K** Kotlin

Introduction
○○○○○○○○○●○○○○
Practical background

Coroutines in details
○○○○○○○○○○○○○○○○○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

# Toy problem wasting computer resources on sleeping

```
1   fun processRequest() {
2       val request = receiveRequest()
3       val processed = handleRequest(request)
4       sendResponse(processed)
5   }
```

**K** Kotlin

Introduction
○○○○○○○○○●○○○○
Practical background

Coroutines in details
○○○○○○○○○○○○○○○○○○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

# Toy problem wasting computer resources on sleeping

```
1   fun processRequest() {
2       val request = receiveRequest()
3       val processed = handleRequest(request)
4       sendResponse(processed)
5   }
```

- potentially, the application may have thousands of these tasks run concurrently

**K** Kotlin

Introduction
○○○○○○○○●○○○○○
Practical background

Coroutines in details
○○○○○○○○○○○○○○○○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

# Toy problem wasting computer resources on sleeping

```
1   fun processRequest() {
2       val request = receiveRequest()
3       val processed = handleRequest(request)
4       sendResponse(processed)
5   }
```

- ■ potentially, the application may have thousands of these tasks run concurrently
- ■ but every of them would consume a single thread to be executed

**K** Kotlin

Introduction
○○○○○○○○○●○○○○
Practical background

Coroutines in details
○○○○○○○○○○○○○○○○○○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

# Toy problem wasting computer resources on sleeping

```
1   fun processRequest() {
2       val request = receiveRequest()
3       val processed = handleRequest(request)
4       sendResponse(processed)
5   }
```

- potentially, the application may have thousands of these tasks run concurrently

- but every of them would consume a single thread to be executed

- the problem are the blocking functions which are related to IO operations and expensive processing

**K** Kotlin

Introduction
○○○○○○○○○○●○○○
Practical background

Coroutines in details
○○○○○○○○○○○○○○○○○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

# The "old" JVM solution for this problem are callbacks

```
1    fun processRequestCallbacks() {
2        receiveRequest { request ->
3            handleRequest(request) { processed ->
4                sendResponse(processed)
5            }
6        }
7    }
```

**K** Kotlin

Introduction
○○○○○○○○○●○○○
Practical background

Coroutines in details
○○○○○○○○○○○○○○○○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

# The "old" JVM solution for this problem are callbacks

```
1  fun processRequestCallbacks() {
2      receiveRequest { request ->
3          handleRequest(request) { processed ->
4              sendResponse(processed)
5          }
6      }
7  }
```

- functions returns immediately and run the callbacks later to get the result

**K** Kotlin

Introduction
○○○○○○○○○○●○○○
Practical background

Coroutines in details
○○○○○○○○○○○○○○○○○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

# The "old" JVM solution for this problem are callbacks

```
1    fun processRequestCallbacks() {
2        receiveRequest { request ->
3            handleRequest(request) { processed ->
4                sendResponse(processed)
5            }
6        }
7    }
```

- functions returns immediately and run the callbacks later to get the result
- this introduces "callback hell" problem

**K** Kotlin

Introduction
○○○○○○○○○○●○○○
Practical background

Coroutines in details
○○○○○○○○○○○○○○○○○○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

# The "old" JVM solution for this problem are callbacks

```
1  fun processRequestCallbacks() {
2      receiveRequest { request ->
3          handleRequest(request) { processed ->
4              sendResponse(processed)
5          }
6      }
7  }
```

- functions returns immediately and run the callbacks later to get the result
- this introduces "callback hell" problem
- this code has multiple limitations - the structure of code prevents imperative coding (how to write loops?)

**K** Kotlin

Introduction
○○○○○○○○○○○●○○

Coroutines in details
○○○○○○○○○○○○○○○○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

Practical background

# Writing imperative code with callbacks

```
1   fun imperative() {
2       var number = 0
3       for (i in 1..1_000)
4           number = inc(number)
5       println(number)
6   }
```

K Kotlin

Introduction
○○○○○○○○○○○●○○
Practical background

Coroutines in details
○○○○○○○○○○○○○○○○○○○○○○

Higher-level APIs
○○○○○○○○○

Appendix
○○

# Writing imperative code with callbacks

```kotlin
 1   fun imperativeCallback() {
 2       class LateInitCallbackWrapper<T : Any> {
 3           lateinit var callback: (T) -> Unit
 4       }
 5
 6       val calls = Array(1_000) { LateInitCallbackWrapper<Int>() }
 7       for (i in 0..998)
 8           calls[i].callback = { incWithCallback(it, calls[i + 1].callback::invoke) }
 9       calls[999].callback = { incWithCallback(it, ::println) }
10       calls[0].callback(0)
11   }
```

**K** Kotlin

**Introduction**
oooooooooo●oo

Coroutines in details
ooooooooooooooooooooooo

Higher-level APIs
oooooooo

Appendix
oo

Practical background

# Writing imperative code with callbacks

```
 1   fun imperativeCallback() {
 2       class LateInitCallbackWrapper<T : Any> {
 3           lateinit var callback: (T) -> Unit
 4       }
 5
 6       val calls = Array(1_000) { LateInitCallbackWrapper<Int>() }
 7       for (i in 0..998)
 8           calls[i].callback = { incWithCallback(it, calls[i + 1].callback::invoke) }
 9       calls[999].callback = { incWithCallback(it, ::println) }
10       calls[0].callback(0)
11   }
```

- hard to understand intention of programmer using callbacks

**K** Kotlin

Introduction      Coroutines in details      Higher-level APIs      Appendix
○○○○○○○○○○●○○    ○○○○○○○○○○○○○○○○○○○○○○    ○○○○○○○○     ○○
Practical background

# Writing imperative code with callbacks

```
 1   fun imperativeCallback() {
 2       class LateInitCallbackWrapper<T : Any> {
 3           lateinit var callback: (T) -> Unit
 4       }
 5
 6       val calls = Array(1_000) { LateInitCallbackWrapper<Int>() }
 7       for (i in 0..998)
 8           calls[i].callback = { incWithCallback(it, calls[i + 1].callback::invoke) }
 9       calls[999].callback = { incWithCallback(it, ::println) }
10       calls[0].callback(0)
11   }
```

- hard to understand intention of programmer using callbacks
- additionally, there's extra cost related to working with the references for callbacks

**K** Kotlin

Introduction
○○○○○○○○○○●○○

Coroutines in details
○○○○○○○○○○○○○○○○○○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

Practical background

# Writing imperative code with callbacks

```kotlin
 1   fun imperativeCallback() {
 2       class LateInitCallbackWrapper<T : Any> {
 3           lateinit var callback: (T) -> Unit
 4       }
 5
 6       val calls = Array(1_000) { LateInitCallbackWrapper<Int>() }
 7       for (i in 0..998)
 8           calls[i].callback = { incWithCallback(it, calls[i + 1].callback::invoke) }
 9       calls[999].callback = { incWithCallback(it, ::println) }
10       calls[0].callback(0)
11   }
```

- hard to understand intention of programmer using callbacks
- additionally, there's extra cost related to working with the references for callbacks
- this solution doesn't scale well so other approach needs to be found

K Kotlin

Introduction
○○○○○○○○○○○○●○
Practical background

Coroutines in details
○○○○○○○○○○○○○○○○○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

# Futures/Promises/Rx being a rescue?

```
1  fun processRequestPromise() {
2      receiveRequestPromise()
3          .flatMap { handleRequestPromise(it) }
4          .subscribeBy(
5              onSuccess = { sendResponse(it) },
6              onError = { println("error: $it") },
7          )
8  }
```

**K** Kotlin

Introduction
○○○○○○○○○○○●○
Practical background

Coroutines in details
○○○○○○○○○○○○○○○○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

# Futures/Promises/Rx being a rescue?

```
1  fun processRequestPromise() {
2      receiveRequestPromise()
3          .flatMap { handleRequestPromise(it) }
4          .subscribeBy(
5              onSuccess = { sendResponse(it) },
6              onError = { println("error: $it") },
7          )
8  }
```

- we got nicer, composable code, with no extra indentation

**K** Kotlin

Introduction
○○○○○○○○○○○○●○

Coroutines in details
○○○○○○○○○○○○○○○○○○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

Practical background

# Futures/Promises/Rx being a rescue?

```
1  fun processRequestPromise() {
2      receiveRequestPromise()
3          .flatMap { handleRequestPromise(it) }
4          .subscribeBy(
5              onSuccess = { sendResponse(it) },
6              onError = { println("error: $it") },
7          )
8  }
```

- we got nicer, composable code, with no extra indentation
- we're able to propagate our exceptions

**K** Kotlin

Introduction
○○○○○○○○○○○○○●○
Coroutines in details
○○○○○○○○○○○○○○○○○○○○○○○
Higher-level APIs
○○○○○○○○
Appendix
○○

Practical background

# Futures/Promises/Rx being a rescue?

```
1   fun processRequestPromise() {
2       receiveRequestPromise()
3           .flatMap { handleRequestPromise(it) }
4           .subscribeBy(
5               onSuccess = { sendResponse(it) },
6               onError = { println("error: $it") },
7           )
8   }
```

- we got nicer, composable code, with no extra indentation
- we're able to propagate our exceptions
- but the combinators are the problem - it might be hard to learn them all

**K** Kotlin

Introduction
○○○○○○○○○○○○●○

Coroutines in details
○○○○○○○○○○○○○○○○○○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

Practical background

# Futures/Promises/Rx being a rescue?

```
1   fun processRequestPromise() {
2       receiveRequestPromise()
3           .flatMap { handleRequestPromise(it) }
4           .subscribeBy(
5               onSuccess = { sendResponse(it) },
6               onError = { println("error: $it") },
7           )
8   }
```

- we got nicer, composable code, with no extra indentation

- we're able to propagate our exceptions

- but the combinators are the problem - it might be hard to learn them all

- while there are different libraries with other names

**K** Kotlin

# Kotlin solution are coroutines

```
1   suspend fun processRequestCoroutines() {
2       val request = receiveRequestCoroutines()
3       val processed = handleRequestCoroutines(request)
4       sendResponseCoroutines(processed)
5   }
```

**K** Kotlin

Introduction
○○○○○○○○○○○○○●
Practical background

Coroutines in details
○○○○○○○○○○○○○○○○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

# Kotlin solution are coroutines

```
1   suspend fun processRequestCoroutines() {
2       val request = receiveRequestCoroutines()
3       val processed = handleRequestCoroutines(request)
4       sendResponseCoroutines(processed)
5   }
```

- need to mark functions with `suspend` modifier

**K** Kotlin

Introduction
○○○○○○○○○○○○●○
Practical background

Coroutines in details
○○○○○○○○○○○○○○○○○○○○○○

Higher-level APIs
○○○○○○○○○

Appendix
○○

# Kotlin solution are coroutines

```
1   suspend fun processRequestCoroutines() {
2       val request = receiveRequestCoroutines()
3       val processed = handleRequestCoroutines(request)
4       sendResponseCoroutines(processed)
5   }
```

- need to mark functions with `suspend` modifier
- but the structure of the code looks like regular code (loops, exceptions, extension functions)

**K** Kotlin

Introduction
○○○○○○○○○○○○○●

Coroutines in details
○○○○○○○○○○○○○○○○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

Practical background

# Kotlin solution are coroutines

```
1    suspend fun processRequestCoroutines() {
2        val request = receiveRequestCoroutines()
3        val processed = handleRequestCoroutines(request)
4        sendResponseCoroutines(processed)
5    }
```

- need to mark functions with `suspend` modifier
- but the structure of the code looks like regular code (loops, exceptions, extension functions)
- we got "not obvious" suspension points

**K** Kotlin

Introduction
○○○○○○○○○○○○●
Practical background

Coroutines in details
○○○○○○○○○○○○○○○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

# Kotlin solution are coroutines

```
1   suspend fun processRequestCoroutines() {
2       val request = receiveRequestCoroutines()
3       val processed = handleRequestCoroutines(request)
4       sendResponseCoroutines(processed)
5   }
```

- need to mark functions with `suspend` modifier
- but the structure of the code looks like regular code (loops, exceptions, extension functions)
- we got "not obvious" suspension points
- but the good IDE helps here a lot if we're interested in this

**K** Kotlin

Introduction
000000000000000

Coroutines in details
●00000000000000000000000

Higher-level APIs
00000000

Appendix
00

# Coroutines in details

Introduction
000000000000

Coroutines in details
0●00000000000000000000

Higher-level APIs
00000000

Appendix
00

Implementation details

# What's under the hood?

- suspending functions are implemented via
  Continuation-Passing-Style (CPS)

**K** Kotlin

Introduction
○○○○○○○○○○○○○○

Coroutines in details
○●○○○○○○○○○○○○○○○○○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

Implementation details

# What's under the hood?

- suspending functions are implemented via Continuation-Passing-Style (CPS)
- every suspending function has an additional `Continuation` parameter passed implicitly when function is invoked

**K** Kotlin

# What's under the hood?

- suspending functions are implemented via Continuation-Passing-Style (CPS)
- every suspending function has an additional `Continuation` parameter passed implicitly when function is invoked
- consider the following signature of `suspend` function

**K** Kotlin

Introduction
000000000000000

Coroutines in details
0●00000000000000000000

Higher-level APIs
00000000

Appendix
00

Implementation details

# What's under the hood?

- suspending functions are implemented via Continuation-Passing-Style (CPS)
- every suspending function has an additional `Continuation` parameter passed implicitly when function is invoked
- consider the following signature of suspend function

```
suspend fun <T> CompletableFuture<T>.await(): T =
```

K Kotlin

Introduction
000000000000

Coroutines in details
0●00000000000000000000

Higher-level APIs
00000000

Appendix
00

Implementation details

# What's under the hood?

- suspending functions are implemented via Continuation-Passing-Style (CPS)
- every suspending function has an additional `Continuation` parameter passed implicitly when function is invoked
- consider the following signature of suspend function

```
suspend fun <T> CompletableFuture<T>.await(): T =
```

- its actual implementation after compiler transformation has the following form

**K** Kotlin

# What's under the hood?

- suspending functions are implemented via Continuation-Passing-Style (CPS)
- every suspending function has an additional `Continuation` parameter passed implicitly when function is invoked
- consider the following signature of suspend function

```
suspend fun <T> CompletableFuture<T>.await(): T =
```

- its actual implementation after compiler transformation has the following form

```
fun <T> CompletableFuture<T>.await(continuation: Continuation<T>): Any? =
```

**K** Kotlin

Introduction
○○○○○○○○○○○○○○○

Coroutines in details
○○●○○○○○○○○○○○○○○○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

Implementation details

# Continuation interface

- `context` represents an arbitrary user-defined context that is associated with the coroutine

```
1   interface Continuation<in T> {
2       val context: CoroutineContext
3       fun resumeWith(result: Result<T>)
4   }
```

**K** Kotlin

# Continuation interface

- `context` represents an arbitrary user-defined context that is associated with the coroutine
- `resumeWith` function is a completion callback that is used to report on coroutine completion:

```
1    interface Continuation<in T> {
2        val context: CoroutineContext
3        fun resumeWith(result: Result<T>)
4    }
```

**K** Kotlin

Introduction
0000000000000

Coroutines in details
0000000000000000000000

Higher-level APIs
00000000

Appendix
00

Implementation details

# Continuation interface

- **context** represents an arbitrary user-defined context that is associated with the coroutine
- **resumeWith** function is a completion callback that is used to report on coroutine completion:
  - a success - with a value

```
1   interface Continuation<in T> {
2       val context: CoroutineContext
3       fun resumeWith(result: Result<T>)
4   }
```

**K** Kotlin

Introduction
0000000000000
Coroutines in details
00●0000000000000000000
Higher-level APIs
00000000
Appendix
00
Implementation details

# Continuation interface

- **context** represents an arbitrary user-defined context that is associated with the coroutine
- **resumeWith** function is a completion callback that is used to report on coroutine completion:
    - a success - with a value
    - a failure - with an exception

```
1    interface Continuation<in T> {
2        val context: CoroutineContext
3        fun resumeWith(result: Result<T>)
4    }
```

**K** Kotlin

Introduction
000000000000

Coroutines in details
00●0000000000000000000

Higher-level APIs
00000000

Appendix
00

Implementation details

# Continuation interface

- **context** represents an arbitrary user-defined context that is associated with the coroutine
- **resumeWith** function is a completion callback that is used to report on coroutine completion:
    - a success - with a value
    - a failure - with an exception
- there are two extension functions defined for convenience

```
1  interface Continuation<in T> {
2      val context: CoroutineContext
3      fun resumeWith(result: Result<T>)
4  }
```

```
1  fun <T> Continuation<T>.resume(value: T) { .. }
2
3  fun <T> Continuation<T>.resumeWithException(exception: Throwable) { .. }
```

**K** Kotlin

Introduction
ooooooooooooo

Coroutines in details
ooooooooooooooooooooo

Higher-level APIs
oooooooo

Appendix
oo

Implementation details

# How the typical implementation looks like?

```
1   suspend fun <T> CompletableFuture<T>.await(): T =
2       suspendCoroutine { cont: Continuation<T> ->
3           whenComplete { result, exception ->
4               if (exception == null) cont.resume(result)
5               else cont.resumeWithException(exception)
6           }
7       }
```

K Kotlin

Introduction
ooooooooooooo
Coroutines in details
ooooOooooooooooooooooo
Higher-level APIs
ooooooooo
Appendix
oo
Implementation details

# How the typical implementation looks like?

```
 1    suspend fun <T> CompletableFuture<T>.await(): T =
 2        suspendCoroutine { cont: Continuation<T> ->
 3            (this as CompletableFuture<T>).whenComplete(
 4                object : BiConsumer<T, Throwable?> {
 5                    override fun accept(result: T, exception: Throwable?) {
 6                        if (exception == null) cont.resume(result)
 7                        else cont.resumeWithException(exception)
 8                    }
 9                }
10            )
11        }
```

K Kotlin

# How the typical implementation looks like?

```
1    suspend fun <T> CompletableFuture<T>.await(): T =
2        suspendCoroutine { cont: Continuation<T> ->
3            (this as CompletableFuture<T>).whenComplete(
4                object : BiConsumer<T, Throwable?> {
5                    override fun accept(result: T, exception: Throwable?) {
6                        if (exception == null) cont.resume(result)
7                        else cont.resumeWithException(exception)
8                    }
9                }
10           )
11       }
```

- we use the callback interface to resume the continuation

**K** Kotlin

Introduction
○○○○○○○○○○○○○

Coroutines in details
○○○●○○○○○○○○○○○○○○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

Implementation details

# How the typical implementation looks like?

```
1    suspend fun <T> CompletableFuture<T>.await(): T =
2        suspendCoroutine { cont: Continuation<T> ->
3            (this as CompletableFuture<T>).whenComplete(
4                object : BiConsumer<T, Throwable?> {
5                    override fun accept(result: T, exception: Throwable?) {
6                        if (exception == null) cont.resume(result)
7                        else cont.resumeWithException(exception)
8                    }
9                }
10            )
11        }
```

- we use the callback interface to resume the continuation
- to actually suspend execution function must invoke other suspending function - `await` implementation invokes a suspending function suspendCoroutine

**K** Kotlin

Introduction
000000000000
Coroutines in details
00000●000000000000000000
Higher-level APIs
00000000
Appendix
00

Implementation details

# State machines as perfect solution

- it's crucial to implement coroutines efficiently, i.e. create as few classes and objects as possible

K Kotlin

# State machines as perfect solution

- it's crucial to implement coroutines efficiently, i.e. create as few classes and objects as possible

- compiler creates a single instance of class that may have any number of suspension points

**K** Kotlin

Introduction
000000000000000
Coroutines in details
0000000000000000000000
Higher-level APIs
00000000
Appendix
00

Implementation details

# State machines as perfect solution

- it's crucial to implement coroutines efficiently, i.e. create as few classes and objects as possible

- compiler creates a single instance of class that may have any number of suspension points

- general idea - suspending function is a state machine

**K** Kotlin

Introduction
0000000000000

Coroutines in details
0000●000000000000000

Higher-level APIs
00000000

Appendix
00

Implementation details

# State machines as perfect solution

- it's crucial to implement coroutines efficiently, i.e. create as few classes and objects as possible
- compiler creates a single instance of class that may have any number of suspension points
- general idea - suspending function is a state machine
- states corresponds to suspension points

**K** Kotlin

Introduction
000000000000000

Coroutines in details
0000●00000000000000000

Higher-level APIs
00000000

Appendix
00

Implementation details

# State machines as perfect solution

- it's crucial to implement coroutines efficiently, i.e. create as few classes and objects as possible

- compiler creates a single instance of class that may have any number of suspension points

- general idea - suspending function is a state machine

- states corresponds to suspension points

```
1   suspend fun stateMachine() {
2       val a = a()
3       val y = foo(a).await() // 1st suspension point
4       b()
5       val z = bar(a, y).await() // 2nd suspension point
6       c(z)
7   }
```

**K** Kotlin

Introduction
000000000000
Coroutines in details
00000●0000000000000000
Higher-level APIs
00000000
Appendix
00
Implementation details

# State machines as perfect solution

- three states for the given block of code:

**K** Kotlin

# State machines as perfect solution

- three states for the given block of code:
  1. initial (before any suspension point)

**K** Kotlin

Introduction
0000000000000

Coroutines in details
0000000000000000000000000

Higher-level APIs
00000000

Appendix
00

Implementation details

# State machines as perfect solution

- three states for the given block of code:
  1. initial (before any suspension point)
  2. after the first suspension point

**K** Kotlin

Introduction
000000000000
Coroutines in details
00000●0000000000000000
Higher-level APIs
00000000
Appendix
00

Implementation details

# State machines as perfect solution

- three states for the given block of code:
  1. initial (before any suspension point)
  2. after the first suspension point
  3. after the second suspension point

**K** Kotlin

# State machines as perfect solution

- three states for the given block of code:
    1. initial (before any suspension point)
    2. after the first suspension point
    3. after the second suspension point

- each of them being an entry point to continuation of this block

**K** Kotlin

# State machine as perfect solution

- three states for the given block of code:
  1. initial (before any suspension point)
  2. after the first suspension point
  3. after the second suspension point
- each of them being an entry point to continuation of this block
- which can be seen as `switch` statement with label to each suspension point

**K** Kotlin

Introduction
000000000000000

Coroutines in details
0000000000000000000000

Higher-level APIs
00000000

Appendix
00

Implementation details

# State machines as perfect solution

- three states for the given block of code:
  1. initial (before any suspension point)
  2. after the first suspension point
  3. after the second suspension point
- each of them being an entry point to continuation of this block
- which can be seen as `switch` statement with label to each suspension point
- the execution is continued with `resumeWith` method of coroutine which selects the current branch

**K** Kotlin

Introduction
000000000000

Coroutines in details
000000●000000000000000

Higher-level APIs
00000000

Appendix
00

Implementation details

# State machines as perfect solution

- three states for the given block of code:
  1. initial (before any suspension point)
  2. after the first suspension point
  3. after the second suspension point
- each of them being an entry point to continuation of this block
- which can be seen as `switch` statement with label to each suspension point
- the execution is continued with `resumeWith` method of coroutine which selects the current branch
- local variables are generated as class fields of anonymous class

**K** Kotlin

Introduction
000000000000

Coroutines in details
0000000●000000000000000

Higher-level APIs
00000000

Appendix
00

Implementation details

# State machine implementation

Pseudo-code corresponding to JVM bytecode created for
suspending function:

```
1   class <anonymous_for_state_machine> extends SuspendLambda<...> {
2       // current state of the state machine
3       int label = 0
4
5       // local variables of the coroutine
6       A a = null
7       Y y = null
8
9       void resumeWith(Object result) {
10          if (label == 0) goto L0
11          if (label == 1) goto L1
12          if (label == 2) goto L2
13          else throw IllegalStateException()
14          // each label code
```

**K** Kotlin

Introduction
0000000000000

Coroutines in details
0000000●000000000000000

Higher-level APIs
00000000

Appendix
00

Implementation details

# State machines implementation

```
15              // local variables and state switch
16          L0:
17              // result is expected to be `null` at this invocation
18              a = a()
19              label = 1
20              result = foo(a).await(this) // 'this' is passed as a continuation
21              if (result == COROUTINE_SUSPENDED) return // return if await had suspended
22          L1:
23              // external code has resumed this coroutine passing the result of .await()
24              y = (Y) result
25              b()
26              label = 2
27              result = bar(a, y).await(this) // 'this' is passed as a continuation
28              if (result == COROUTINE_SUSPENDED) return // return if await had suspended
29          L2:
30              // external code has resumed this coroutine passing the result of .await()
31              Z z = (Z) result
32              c(z)
33              label = -1 // No more steps are allowed
34              return
35      }
36  }
```

**K** Kotlin

Introduction
○○○○○○○○○○○○○○
Coroutines in details
○○○○○○○○○●○○○○○○○○○○○○○
Higher-level APIs
○○○○○○○○
Appendix
○○

Implementation details

# Generators

- generators with `yield` can be easily implemented using suspend functions

# Generators

- generators with `yield` can be easily implemented using suspend functions
- and that's how the Kotlin std-lib does, providing `Sequence<T>` abstraction with `yield` as a function of `SequenceScope<T>`

**K** Kotlin

# Generators

- generators with `yield` can be easily implemented using suspend functions
- and that's how the Kotlin std-lib does, providing `Sequence<T>` abstraction with `yield` as a function of `SequenceScope<T>`
- they can represent potentially infinite lists (like in Haskell)

**K** Kotlin

Introduction
000000000000
Coroutines in details
000000000●000000000000
Higher-level APIs
00000000
Appendix
00

Implementation details

# Generators

- generators with `yield` can be easily implemented using suspend functions
- and that's how the Kotlin std-lib does, providing `Sequence<T>` abstraction with `yield` as a function of `SequenceScope<T>`
- they can represent potentially infinite lists (like in Haskell)
- or just in general lazily computed sequences

**K** Kotlin

# Generators

- generators with `yield` can be easily implemented using suspend functions
- and that's how the Kotlin std-lib does, providing `Sequence<T>` abstraction with `yield` as a function of `SequenceScope<T>`
- they can represent potentially infinite lists (like in Haskell)
- or just in general lazily computed sequences
- the strength of them is supporting arbitrary control flow, with `for`, `when`, `try/catch` etc.

**K** Kotlin

Introduction
0000000000000

Coroutines in details
000000000●00000000000

Higher-level APIs
00000000

Appendix
00

Implementation details

# Generators - example Fibonacci sequence

```
 1   fun fibonacci() = sequence {
 2       yield(1)
 3       var curr = 1
 4       var next = 1
 5       while (true) {
 6           yield(next)
 7           val tmp = curr + next
 8           curr = next
 9           next = tmp
10       }
11   }
```

**K** Kotlin

Introduction
000000000000000

Coroutines in details
00000000000000000000000000

Higher-level APIs
00000000

Appendix
00

Implementation details

# Generators - example Fibonacci sequence

```
1   fun fibonacci() = sequence {
2       yield(1)
3       var curr = 1
4       var next = 1
5       while (true) {
6           yield(next)
7           val tmp = curr + next
8           curr = next
9           next = tmp
10      }
11  }

    fun main() = fibonacci().take(42).forEach(::println)
```

K Kotlin

Introduction
0000000000000

Coroutines in details
00000000000●00000000000

Higher-level APIs
00000000

Appendix
00

Fundamental concepts

# Initial idea - extra DSL

```
1   fun processRequest() = async {
2       val request = await(receiveRequestPromise())
3       val processed = await(handleRequestPromise(request))
4       sendResponse(processed)
5   }
```

**K** Kotlin

# Initial idea - extra DSL

```
1    fun processRequest() = async {
2        val request = await(receiveRequestPromise())
3        val processed = await(handleRequestPromise(request))
4        sendResponse(processed)
5    }
```

- original prototype defined the library functions, similar to keywords from e.g. C#, Typescript, Dart, Python

**K** Kotlin

Introduction
000000000000

Coutines in details
00000000000●00000000000

Higher-level APIs
00000000

Appendix
00

Fundamental concepts

# Initial idea - extra DSL

```
1   fun processRequest() = async {
2       val request = await(receiveRequestPromise())
3       val processed = await(handleRequestPromise(request))
4       sendResponse(processed)
5   }
```

- original prototype defined the library functions, similar to keywords from e.g. C#, Typescript, Dart, Python
- but the compiler needed to know that await can be suspended, so it was marked with suspend modifier

**K** Kotlin

Introduction
○○○○○○○○○○○○○

Coroutines in details
○○○○○○○○○○○●○○○○○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

Fundamental concepts

# Initial idea - extra DSL

```
1   fun processRequest() = async {
2       val request = await(receiveRequestPromise())
3       val processed = await(handleRequestPromise(request))
4       sendResponse(processed)
5   }
```

- original prototype defined the library functions, similar to keywords from e.g. C#, Typescript, Dart, Python
- but the compiler needed to know that `await` can be suspended, so it was marked with `suspend` modifier
- so the decision was to remove the need for `await` call

**K** Kotlin

Introduction
0000000000000

Coroutines in details
00000000000●00000000000

Higher-level APIs
00000000

Appendix
00

Fundamental concepts

# Initial idea - extra DSL

```
1   fun processRequest() = async {
2       val request = await(receiveRequestPromise())
3       val processed = await(handleRequestPromise(request))
4       sendResponse(processed)
5   }
```

- original prototype defined the library functions, similar to keywords from e.g. C#, Typescript, Dart, Python
- but the compiler needed to know that await can be suspended, so it was marked with suspend modifier
- so the decision was to remove the need for await call
- which led to adding suspend modifier to functions instead of using async builder

**K** Kotlin

Introduction
○○○○○○○○○○○○○
Coroutines in details
○○○○○○○○○○○○●○○○○○○○○○○○
Higher-level APIs
○○○○○○○○
Appendix
○○

Fundamental concepts

# Initial idea - extra DSL

- there was some fear as no other language resigned from
  `async`/`await`

**K** Kotlin

# Initial idea - extra DSL

- there was some fear as no other language resigned from `async`/`await`
- actually in Go language nobody is worried about the not visible suspension points

**K** Kotlin

# Initial idea - extra DSL

- there was some fear as no other language resigned from `async`/`await`
- actually in Go language nobody is worried about the not visible suspension points

```
1   suspend fun go() = coroutineScope {
2       go { printDelayed("go: fizz") }
3       printDelayed("go: buzz")
4   }
```

**K** Kotlin

# Initial idea - extra DSL

- there was some fear as no other language resigned from `async`/`await`
- actually in Go language nobody is worried about the not visible suspension points

```
1   suspend fun go() = coroutineScope {
2       go { printDelayed("go: fizz") }
3       printDelayed("go: buzz")
4   }
```

```
1   suspend fun kt() = coroutineScope {
2       launch { printDelayed("kt: fizz") }
3       printDelayed("kt: buzz")
4   }
```

**K** Kotlin

Introduction                    Coroutines in details              Higher-level APIs        Appendix
oooooooooooooo                  ooooooooooooo●ooooooooooo           ooooooooo                oo
Fundamental concepts

# Initial idea - extra DSL

- there was some fear as no other language resigned from `async`/`await`
- actually in Go language nobody is worried about the not visible suspension points

```
1  suspend fun go() = coroutineScope {
2      go { printDelayed("go: fizz") }
3      printDelayed("go: buzz")
4  }
```

```
1  suspend fun kt() = coroutineScope {
2      launch { printDelayed("kt: fizz") }
3      printDelayed("kt: buzz")
4  }
```

```
1  suspend fun printDelayed(msg: Any, times: Int = 5) {
2      for (i in 1..times) {
3          delay(1_000)
4          println(msg)
5      }
6  }
```

**K** Kotlin

Introduction
○○○○○○○○○○○○○

Coroutines in details
○○○○○○○○○○○○○○●○○○○○○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

Fundamental concepts

# Context and Scope of execution

```
1   suspend fun implicitScope(): Unit = coroutineScope {
2       launch(context = Default) { say("hello") }
3       launch(context = IO) { say("world") }
4   }
```

Introduction                    Coroutines in details              Higher-level APIs              Appendix
oooooooooooooo                  ooooooooooooo**oo**ooooooooo        ooooooooo                     oo
Fundamental concepts

# Context and Scope of execution

```
1    suspend fun implicitScope(): Unit = coroutineScope {
2        launch(context = Default) { say("hello") }
3        launch(context = IO) { say("world") }
4    }
```

- coroutines execute in some context represented by the
  CoroutineContext type

**K** Kotlin

Introduction
○○○○○○○○○○○○○
Coroutines in details
○○○○○○○○○○○○○●○○○○○○○○○○
Higher-level APIs
○○○○○○○○
Appendix
○○

Fundamental concepts

# Context and Scope of execution

```kotlin
1    suspend fun implicitScope(): Unit = coroutineScope {
2        launch(context = Default) { say("hello") }
3        launch(context = IO) { say("world") }
4    }
```

- coroutines execute in some context represented by the `CoroutineContext` type
- it contains the `CoroutineDispatcher` which delegates the coroutine to some thread of execution

**K** Kotlin

Introduction
000000000000

Coroutines in details
00000000000**000000000000

Higher-level APIs
00000000

Appendix
00

Fundamental concepts

# Context and Scope of execution

```
1    suspend fun implicitScope(): Unit = coroutineScope {
2        launch(context = Default) { say("hello") }
3        launch(context = IO) { say("world") }
4    }
```

- coroutines execute in some context represented by the `CoroutineContext` type

- it contains the `CoroutineDispatcher` which delegates the coroutine to some thread of execution

- it's inherited from the parent `CoroutineScope` by default

**K** Kotlin

# Context and Scope of execution

```
1   suspend fun implicitScope(): Unit = coroutineScope {
2       launch(context = Default) { say("hello") }
3       launch(context = IO) { say("world") }
4   }
```

- coroutines execute in some context represented by the
  CoroutineContext type
- it contains the CoroutineDispatcher which delegates the
  coroutine to some thread of execution
- it's inherited from the parent CoroutineScope by default
- and all coroutine builders like launch and async accept
  optional CoroutineContext to explicitly specify the
  dispatcher for the new coroutine context

**K** Kotlin

# Context and Scope of execution

```
1   suspend fun explicitScope(): Unit = coroutineScope {
2       val scope: CoroutineScope = this
3       scope.launch(context = Default) { say("hello") }
4       scope.launch(context = IO) { say("world") }
5   }
```

- coroutines execute in some context represented by the CoroutineContext type
- it contains the CoroutineDispatcher which delegates the coroutine to some thread of execution
- it's inherited from the parent CoroutineScope by default
- and all coroutine builders like launch and async accept optional CoroutineContext to explicitly specify the dispatcher for the new coroutine context

**K** Kotlin

Introduction
○○○○○○○○○○○○○

Coroutines in details
○○○○○○○○○○○○○○●○○○○○○○○

Higher-level APIs
○○○○○○○○○

Appendix
○○

Structured concurrency

# What is structured concurrency?

- compare understanding concurrency to
  understanding program control flow

**K** Kotlin

Introduction
○○○○○○○○○○○○○

Coroutines in details
○○○○○○○○○○○○○●○○○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

Structured concurrency

# What is structured concurrency?

- compare understanding concurrency to understanding program control flow
- with `goto` instruction it can start being hard to understand

sequential    goto

source:
https://vorpus.org/blog/notes-on-structured-concurrency-or-go-statement-considered-harmful/

**K** Kotlin

Introduction
○○○○○○○○○○○○○○

Coroutines in details
○○○○○○○○○○○○○○●○○○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

Structured concurrency

# What is structured concurrency?

- compare understanding concurrency to understanding program control flow

- with `goto` instruction it can start being hard to understand

- but when we limit ourselves to known structures, it's much more obvious and not that limited



source:
https://vorpus.org/blog/notes-on-structured-concurrency-or-go-statement-considered-harmful/

**K** Kotlin

Introduction
○○○○○○○○○○○○○○

Coroutines in details
○○○○○○○○○○○○○○●○○○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

Structured concurrency

# What is structured concurrency?

- compare understanding concurrency to understanding program control flow
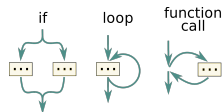- with `goto` instruction it can start being hard to understand
- but when we limit ourselves to known structures, it's much more obvious and not that limited
- in concurrent programming model we have the same problem with tasks

**K** Kotlin

Introduction
○○○○○○○○○○○○○○

Coroutines in details
○○○○○○○○○○○○○●○○○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

Structured concurrency

# What is structured concurrency?

- compare understanding concurrency to understanding program control flow
- with `goto` instruction it can start being hard to understand
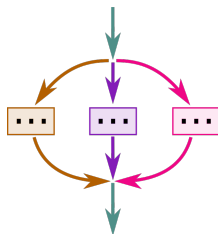- but when we limit ourselves to known structures, it's much more obvious and not that limited
- in concurrent programming model we have the same problem with tasks
- introducing some "structures" to keep everything in consistent state helps in understanding what may happen



source:
https://vorpus.org/blog/notes-on-structured-concurrency-or-go-statement-considered-harmful/

**K** Kotlin

Introduction
○○○○○○○○○○○○○

Coroutines in details
○○○○○○○○○○○○○○●○○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

Structured concurrency

# `launch` and `async` coroutine builders

- `launch` produces a `Job`, while `async` produces a
  `Deferred<T>`

**K** Kotlin

Introduction
000000000000

Coroutines in details
00000000000000●0000000

Higher-level APIs
00000000

Appendix
00

Structured concurrency

# launch and async coroutine builders

- `launch` produces a `Job`, while `async` produces a `Deferred<T>`
- they start a new coroutine without blocking the current one

**K** Kotlin

Introduction
○○○○○○○○○○○○○

Coroutines in details
○○○○○○○○○○○○○○●○○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

Structured concurrency

# launch and async coroutine builders

- launch produces a Job, while async produces a Deferred<T>
- they start a new coroutine without blocking the current one
- they can be called only in context of some CoroutineScope which is "responsible" for its execution

**K** Kotlin

Introduction
○○○○○○○○○○○○○

Coroutines in details
○○○○○○○○○○○○○○●○○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

Structured concurrency

# launch and async coroutine builders

- `launch` produces a `Job`, while `async` produces a `Deferred<T>`
- they start a new coroutine without blocking the current one
- they can be called only in context of some `CoroutineScope` which is "responsible" for its execution
- coroutine is cancelled when the returned `Job` is cancelled

**K** Kotlin

Introduction
0000000000000

Coroutines in details
0000000000000●0000000

Higher-level APIs
00000000

Appendix
00

Structured concurrency

# launch and async coroutine builders

- launch produces a Job, while async produces a Deferred<T>
- they start a new coroutine without blocking the current one
- they can be called only in context of some CoroutineScope which is "responsible" for its execution
- coroutine is cancelled when the returned Job is cancelled
- but still we need to manually check for cancellation in long-running tasks

**K** Kotlin

Introduction
○○○○○○○○○○○○○○

Coroutines in details
○○○○○○○○○○○○○○○○**○○●○○○○○○**

Higher-level APIs
○○○○○○○○

Appendix
○○

Structured concurrency

# Structured concurrency in Kotlin

- `CoroutineScope` can be seen a structure in which the control flows of concurrent tasks are merged

**K** Kotlin

Introduction
○○○○○○○○○○○○○○

Coroutines in details
○○○○○○○○○○○○○○○○●○○○○○○

Higher-level APIs
○○○○○○○○○

Appendix
○○

Structured concurrency

# Structured concurrency in Kotlin

- `CoroutineScope` can be seen a structure in which the control flows of concurrent tasks are merged
- `launch` and `async` corresponds to creating new execution branches

**K** Kotlin

Introduction
○○○○○○○○○○○○○○

Coroutines in details
○○○○○○○○○○○○○○○●○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

Structured concurrency

# Structured concurrency in Kotlin

- `CoroutineScope` can be seen a structure in which the control flows of concurrent tasks are merged
- `launch` and `async` corresponds to creating new execution branches
- parent always waits for children completion

**K** Kotlin

Introduction
○○○○○○○○○○○○○

Coroutines in details
○○○○○○○○○○○○○○○●○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

Structured concurrency

# Structured concurrency in Kotlin

- `CoroutineScope` can be seen a structure in which the control flows of concurrent tasks are merged
- `launch` and `async` corresponds to creating new execution branches
- parent always waits for children completion
- there's no place for loosing some resource

**K** Kotlin

# Structured concurrency in Kotlin

- `CoroutineScope` can be seen a structure in which the control flows of concurrent tasks are merged
- `launch` and `async` corresponds to creating new execution branches
- parent always waits for children completion
- there's no place for loosing some resource
- there's no exception lost - they're propagated

**K** Kotlin

Introduction
○○○○○○○○○○○○○

Coroutines in details
○○○○○○○○○○○○○○○●○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

Structured concurrency

# launch running examples

```
1   suspend fun justSuspensionPoints(): Unit = coroutineScope {
2       sayA()
3       sayB()
4   }
```

**K** Kotlin

Introduction
○○○○○○○○○○○○○○

Coroutines in details
○○○○○○○○○○○○○○○●○○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

Structured concurrency

# `launch` running examples

```
1   suspend fun justSuspensionPoints(): Unit = coroutineScope {
2       sayA()
3       sayB()
4   }
```

- executes sayA and sayB sequentially

**K** Kotlin

Introduction
○○○○○○○○○○○○○

Coroutines in details
○○○○○○○○○○○○○○○●○○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

Structured concurrency

# `launch` running examples

```
1   suspend fun justSuspensionPoints(): Unit = coroutineScope {
2       sayA()
3       sayB()
4   }
```

- executes `sayA` and `sayB` sequentially
- finishes `justSuspensionPoints` when both jobs are finished

**K** Kotlin

Introduction
○○○○○○○○○○○○○

Coroutines in details
○○○○○○○○○○○○○○○○○○●○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

Structured concurrency

# `launch` running examples

```kotlin
1  suspend fun launchFirstInScope(): Unit = coroutineScope {
2      launch { sayA() }
3      sayB()
4  }
```

**K** Kotlin

Introduction
○○○○○○○○○○○○○

Coroutines in details
○○○○○○○○○○○○○○○○○●○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

Structured concurrency

# `launch` running examples

```
1   suspend fun launchFirstInScope(): Unit = coroutineScope {
2       launch { sayA() }
3       sayB()
4   }
```

- executes sayA and sayB concurrently

**K** Kotlin

Introduction
○○○○○○○○○○○○○○

Coroutines in details
○○○○○○○○○○○○○○○○○○○●○○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

Structured concurrency

# `launch` running examples

```
1    suspend fun launchFirstInScope(): Unit = coroutineScope {
2        launch { sayA() }
3        sayB()
4    }
```

- executes `sayA` and `sayB` concurrently
- finishes `launchFirstInScope` when both jobs are finished

**K** Kotlin

Introduction
oooooooooooooo

Coroutines in details
ooooooooooooooo●ooo

Higher-level APIs
oooooooo

Appendix
oo

Structured concurrency

# launch running examples

```
1  suspend fun launchBothInScope(): Unit = coroutineScope {
2      launch { sayA() }
3      launch { sayB() }
4  }
```

**K** Kotlin

Introduction
○○○○○○○○○○○○○

Coroutines in details
○○○○○○○○○○○○○○○○○●○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

Structured concurrency

# `launch` running examples

```
1    suspend fun launchBothInScope(): Unit = coroutineScope {
2        launch { sayA() }
3        launch { sayB() }
4    }
```

- executes `sayA` and `sayB` concurrently

**K** Kotlin

Introduction
○○○○○○○○○○○○○○

Coroutines in details
○○○○○○○○○○○○○○○○○○○○○●○○○

Higher-level APIs
○○○○○○○○

Appendix
○○

Structured concurrency

# `launch` running examples

```
1   suspend fun launchBothInScope(): Unit = coroutineScope {
2       launch { sayA() }
3       launch { sayB() }
4   }
```

- executes `sayA` and `sayB` concurrently
- finishes `launchBothInScope` when both jobs are finished

**K** Kotlin

Introduction
○○○○○○○○○○○○○○

Coroutines in details
○○○○○○○○○○○○○○○○●○○

Higher-level APIs
○○○○○○○○

Appendix
○○

Structured concurrency

# launch running examples

```kotlin
1   suspend fun launchWithException(): Unit = coroutineScope {
2       launch { error("illegal to sayA") }
3       launch { sayB() }
4   }
```

**K** Kotlin

Introduction
○○○○○○○○○○○○○

Coroutines in details
○○○○○○○○○○○○○○○●○○

Higher-level APIs
○○○○○○○○

Appendix
○○

Structured concurrency

# `launch` running examples

```
1   suspend fun launchWithException(): Unit = coroutineScope {
2       launch { error("illegal to sayA") }
3       launch { sayB() }
4   }
```

- executes `sayA` and `sayB` concurrently

**K** Kotlin

Introduction
○○○○○○○○○○○○○

Coroutines in details
○○○○○○○○○○○○○○○○○○○○●○○

Higher-level APIs
○○○○○○○○

Appendix
○○

Structured concurrency

# `launch` running examples

```
1   suspend fun launchWithException(): Unit = coroutineScope {
2       launch { error("illegal to sayA") }
3       launch { sayB() }
4   }
```

- executes sayA and sayB concurrently

- when sayA fails, sayB is cancelled as well

**K** Kotlin

Introduction
○○○○○○○○○○○○○○

Coroutines in details
○○○○○○○○○○○○○○○●○

Higher-level APIs
○○○○○○○○○

Appendix
○○

Structured concurrency

# launch running examples

```
1  suspend fun launchInSupervisorWithException(): Unit = supervisorScope {
2      launch { error("illegal to sayA") }
3      launch { sayB() }
4  }
```

**K** Kotlin

Introduction
0000000000000

Coroutines in details
00000000000000000000000●0

Higher-level APIs
00000000

Appendix
00

Structured concurrency

# `launch` running examples

```
1  suspend fun launchInSupervisorWithException(): Unit = supervisorScope {
2      launch { error("illegal to sayA") }
3      launch { sayB() }
4  }
```

- executes `sayA` and `sayB` concurrently

**K** Kotlin

Introduction
○○○○○○○○○○○○○○

Coroutines in details
○○○○○○○○○○○○○○○●○○○○○○○●○

Higher-level APIs
○○○○○○○○

Appendix
○○

Structured concurrency

# `launch` running examples

```
1   suspend fun launchInSupervisorWithException(): Unit = supervisorScope {
2       launch { error("illegal to sayA") }
3       launch { sayB() }
4   }
```

- executes `sayA` and `sayB` concurrently
- when `sayA` fails, `sayB` is executed until it finishes

**K** Kotlin

Introduction
000000000000

Coroutines in details
000000000000000000000000●

Higher-level APIs
00000000

Appendix
00

Structured concurrency

# Higher-order function

- introducing coroutine brings advantage of no extra need for higher-order combinators like `flatMap`

**K** Kotlin

Introduction
○○○○○○○○○○○○○

Coroutines in details
○○○○○○○○○○○○○○○○○○○○●

Higher-level APIs
○○○○○○○○

Appendix
○○

Structured concurrency

# Higher-order function

- introducing coroutine brings advantage of no extra need for higher-order combinators like `flatMap`
- it's enough to use the `suspend` modifier for lambda parameter

**K** Kotlin

Introduction
oooooooooooooo

Coroutines in details
ooooooooooooooo●oooooooo●

Higher-level APIs
ooooooooo

Appendix
oo

Structured concurrency

# Higher-order function

- introducing coroutine brings advantage of no extra need for higher-order combinators like `flatMap`
- it's enough to use the `suspend` modifier for lambda parameter

```
1   suspend fun main() = retry(afterMillis = 10) { say("hello") }
2
3   suspend fun retry(afterMillis: Long, action: suspend () -> Unit) {
4       do {
5           try {
6               return action()
7           } catch (e: Exception) {
8               println(e.stackTrace)
9           }
10          delay(afterMillis)
11      } while (coroutineContext.isActive)
12  }
```

**K** Kotlin

Introduction
○○○○○○○○○○○○○○

Coroutines in details
○○○○○○○○○○○○○○○○○○○○○○○○

Higher-level APIs
●○○○○○○○

Appendix
○○

# Higher-level APIs

# Channel<T> in Kotlin

- is implemented in Kotlin as a library, conceptually very similar to BlockingQueue

**K** Kotlin

# Channel<T> in Kotlin

- is implemented in Kotlin as a library, conceptually very similar to `BlockingQueue`

- implements `SendChannel<T>` and `ReceiveChannel<T>`

```
1    interface SendChannel<T> {
2        suspend fun send(value: T)
3        fun close()
4    }
5
6    interface ReceiveChannel<T> {
7        suspend fun receive(): T
8        suspend operator fun iterator(): ReceiveIterator<T>
9    }
```

**K** Kotlin

# Channel<T> in Kotlin

- is implemented in Kotlin as a library, conceptually very similar to BlockingQueue
- implements SendChannel<T> and ReceiveChannel<T>
- send suspends when the channel buffer is full, while receive suspends when the buffer is empty

**K** Kotlin

# Channel<T> in Kotlin

- is implemented in Kotlin as a library, conceptually very similar to `BlockingQueue`
- implements `SendChannel<T>` and `ReceiveChannel<T>`
- `send` suspends when the channel buffer is full, while `receive` suspends when the buffer is empty
- channels are hot - there's a coroutine on the other side of the channel that produces values, so we cannot just drop a reference to the ReceiveChannel, because the producer is going to be suspended forever waiting for a consumer, wasting memory resources, open network connections, etc.

**K** Kotlin

Introduction
○○○○○○○○○○○○○○

Coroutines in details
○○○○○○○○○○○○○○○○○○○○○○○○

Higher-level APIs
○○●○○○○○

Appendix
○○

Main ideas

# Flow<T> in Kotlin

- represents the stream of values that are being computed asynchronously

**K** Kotlin

Introduction
○○○○○○○○○○○○○
Coroutines in details
○○○○○○○○○○○○○○○○○○○○○○○
**Higher-level APIs**
○○●○○○○○○
Appendix
○○

Main ideas

# Flow<T> in Kotlin

- represents the stream of values that are being computed asynchronously
- can be used just like a Sequence<T> type for synchronously computed values

**K** Kotlin

Introduction
००००००००००००००

Coroutines in details
००००००००००००००००००००००००

Higher-level APIs
००●००००००

Appendix
००

Main ideas

# Flow<T> in Kotlin

- represents the stream of values that are being computed asynchronously
- can be used just like a Sequence$<T>$ type for synchronously computed values
- flows are cold streams similar to sequences

K Kotlin

Introduction
0000000000000

Coroutines in details
00000000000000000000000

**Higher-level APIs**
0000000000

Appendix
00

Main ideas

# Flow<T> in Kotlin

- represents the stream of values that are being computed asynchronously
- can be used just like a Sequence<T> type for synchronously computed values
- flows are cold streams similar to sequences
- the code inside a flow builder does not run until the flow is collected

**K** Kotlin

Introduction
○○○○○○○○○○○○○○

Coroutines in details
○○○○○○○○○○○○○○○○○○○○○○

Higher-level APIs
○○○●○○○○

Appendix
○○

Practical example

# Images downloader - goals

- have single abstraction responsible for caching the results

# Images downloader - goals

- have single abstraction responsible for caching the results
- have multiple workers responsible for doing the job

**K** Kotlin

Introduction
00000000000000

Coroutines in details
0000000000000000000000

Higher-level APIs
0000●0000

Appendix
00

Practical example

# Images downloader - goals

- have single abstraction responsible for caching the results
- have multiple workers responsible for doing the job
- communicate between these abstractions in safely way

**K** Kotlin

Introduction
○○○○○○○○○○○○○○

Coroutines in details
○○○○○○○○○○○○○○○○○○○○○○○

Higher-level APIs
○○○●○○○○○

Appendix
○○

Practical example

# Images downloader - goals

- have single abstraction responsible for caching the results
- have multiple workers responsible for doing the job
- communicate between these abstractions in safely way
- structure our concurrency model to be sure what's going on

**K** Kotlin

# Images downloader - downloader

```kotlin
 1   fun CoroutineScope.downloader(
 2       receiveStringUrlChannel: ReceiveChannel<String>,
 3       sendResultByteArrayChannel: SendChannel<ByteArray>,
 4       sendToWorkerChannel: SendChannel<Url>,
 5       receiveFromWorkerChannel: ReceiveChannel<DownloadedData>,
 6   ) = launch {
 7       val cached = mutableMapOf<Url, ByteArray>()
 8       while (isActive) {
 9           select {
10               receiveFromWorkerChannel.onReceive { (url, data) ->
11                   cached[url] = data
12                   sendResultByteArrayChannel.send(data)
13               }
14               receiveStringUrlChannel.onReceive {
15                   val url = Url(it)
16                   val data = cached[url]
17                   if (data == null) sendToWorkerChannel.send(url)
18                   else sendResultByteArrayChannel.send(data)
19               }
20           }
21       }
22   }
```

**K** Kotlin

Introduction
○○○○○○○○○○○○○

Coroutines in details
○○○○○○○○○○○○○○○○○○○○○○○

Higher-level APIs
○○○○○○●●○○

Appendix
○○

Practical example

# Images downloader - worker

```
1   fun CoroutineScope.downloadWorker(
2       receive: ReceiveChannel<Url>,
3       sendData: SendChannel<DownloadedData>,
4   ) = launch {
5       val client = HttpClient(CIO)
6       for (url in receive) {
7           val data = client.get(url)
8           val bytes = data.readBytes()
9           val downloaded = DownloadedData(url, bytes)
10          sendData.send(downloaded)
11      }
12  }
```

**K** Kotlin

# Images downloader - common scope
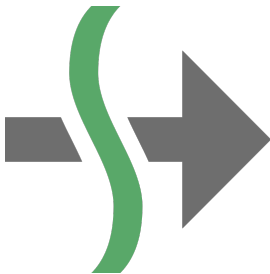
```
1   fun CoroutineScope.processUrls(
2       receiveStringUrlChannel: ReceiveChannel<String>,
3       sendResultByteArrayChannel: SendChannel<ByteArray>,
4   ) {
5       val urls = Channel<Url>(capacity = 1)
6       val data = Channel<DownloadedData>(capacity = 1)
7       repeat(N_WORKERS) { downloadWorker(urls, data) }
8       downloader(receiveStringUrlChannel, sendResultByteArrayChannel, urls, data)
9   }
```

**K** Kotlin

Introduction
○○○○○○○○○○○○○
Coroutines in details
○○○○○○○○○○○○○○○○○○○○○○○
**Higher-level APIs**
○○○○○○○○●
Appendix
○○

Practical example

# Images downloader - Jetpack Compose demo

```
1   fun main() = application {
2       val scope = rememberCoroutineScope()
3       val urlChannel = Channel<String>(capacity = 1)
4       val imageDataChannel = Channel<ByteArray>(capacity = 1)
5       LaunchedEffect(Unit) { processUrls(urlChannel, imageDataChannel) }
6       // application ui
```

**K** Kotlin

# Still confused? Check out these places

1. Roman Elizarov - Structured concurrency  ▸ Watch
2. KotlinConf 2017 - Introduction to Coroutines by Roman Elizarov  ▸ Watch
3. KotlinConf 2017 - Deep Dive into Coroutines on JVM by Roman Elizarov  ▸ Watch
4. KotlinConf 2018 - Kotlin Coroutines in Practice by Roman Elizarov  ▸ Watch
5. KotlinConf 2019 - Asynchronous Data Streams with Kotlin Flow by Roman Elizarov  ▸ Watch
6. Notes on structured concurrency, or: Go statement considered harmful  ▸ Read
7. Revisiting Coroutines, Ana Moura & Roberto Ierusalimschy  ▸ Read

K Kotlin

*Time to suspend*

*Thank you for your attention!*