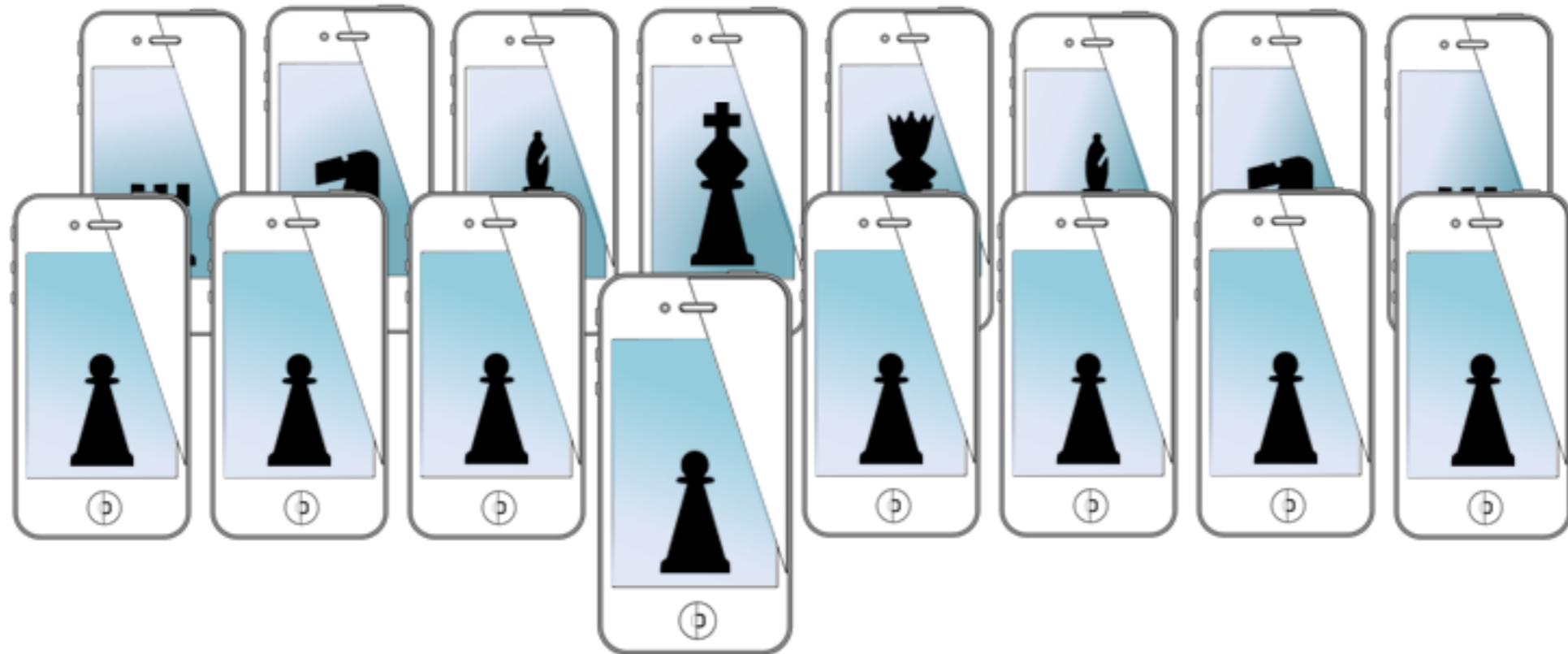


MOBILE SENSING LEARNING



CS5323 & 7323
Mobile Sensing and Learning

More Advanced ARKit and Sound Classification

Eric C. Larson, Lyle School of Engineering,
Computer Science, Southern Methodist University

course logistics

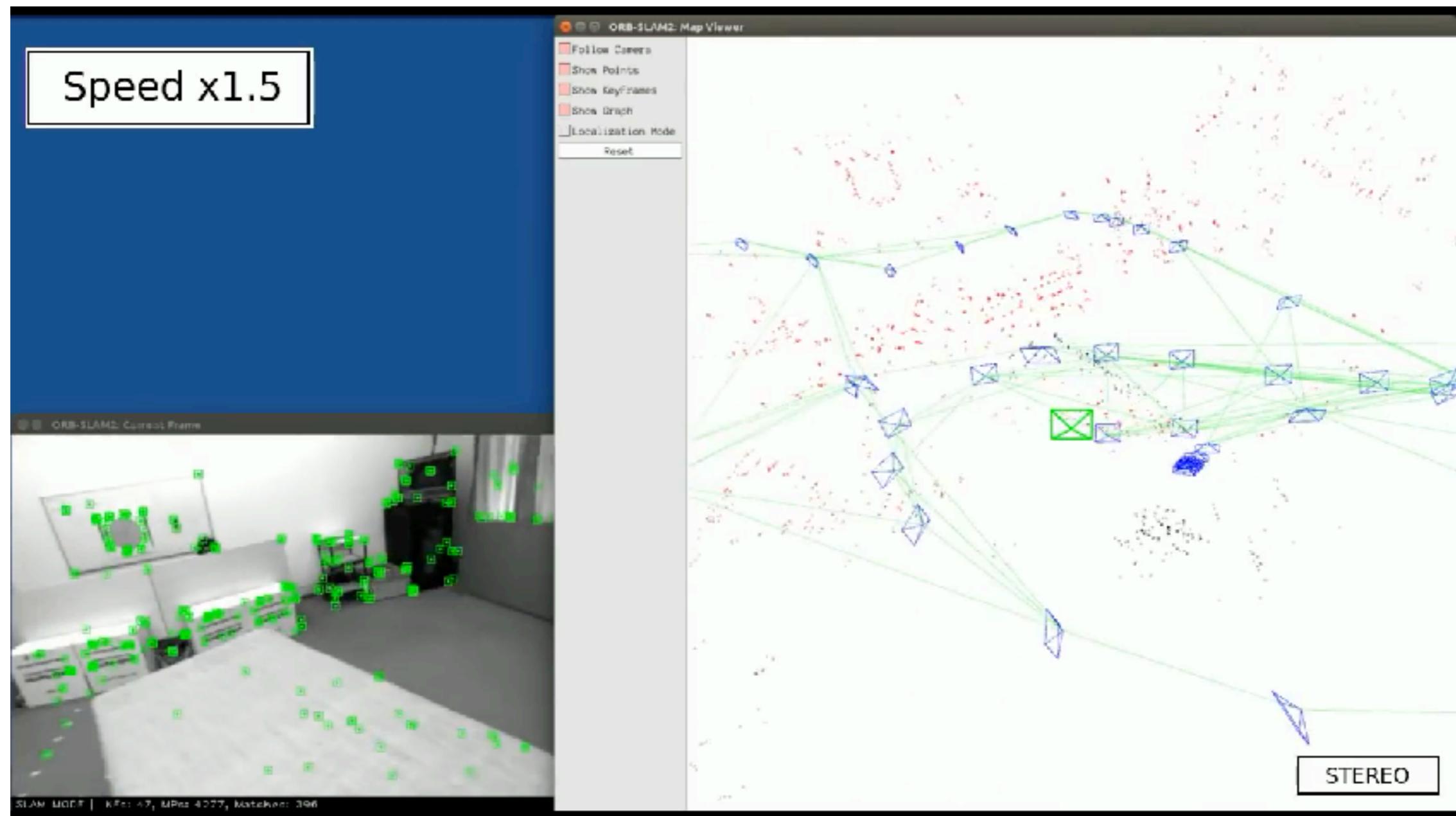
- ARKit 2D and 3D Object Recognition
- More advanced ARKit UI Elements
 - Text
 - Objects
 - Animation
 - Video
- Sound Classification

Why perform detection?

- need to perform detection of images and 3D objects
 - to anchor AR in the world
 - for shared AR experiences
- keep in mind
 - automatic object detection is a slow process
 - ...but is very easy to get started

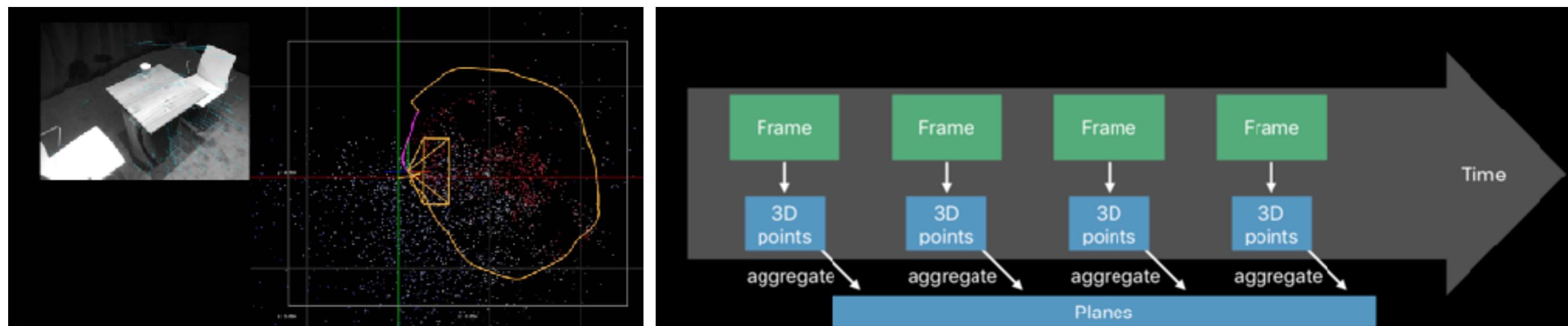
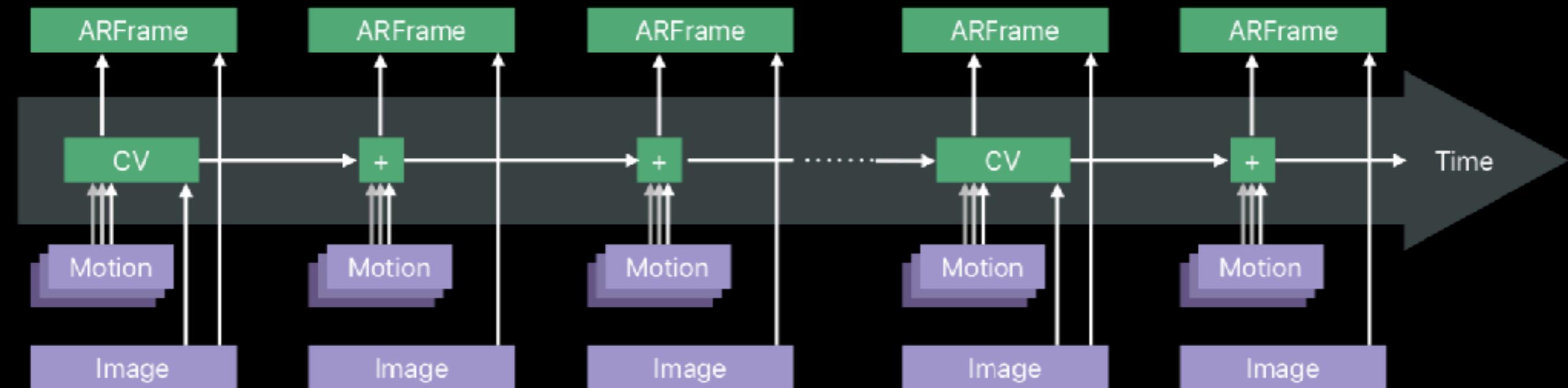
SLAM

- ORB-SLAM: Uses Key points in Image
- Simultaneous Localization and Mapping



SLAM in iOS: World Map

Motion data and computer vision



<https://developer.apple.com/videos/play/wwdc2018/610/>

image recognition in ARKit

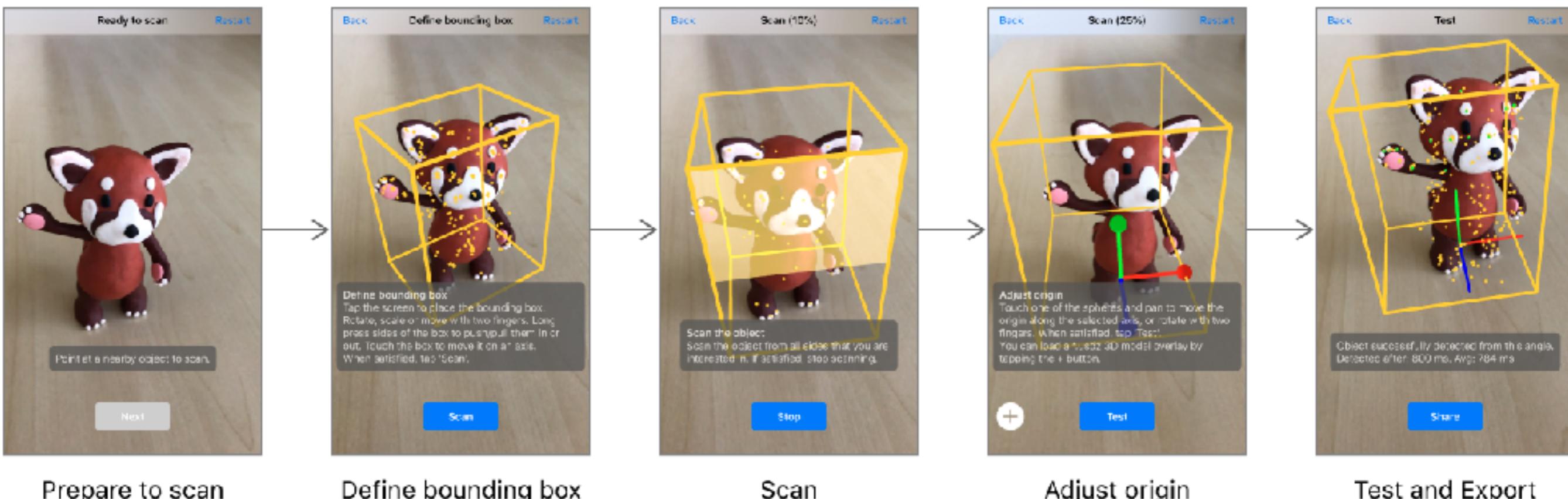
- detection based on scans of keypoints in example image, projected onto detected planes
- drag and drop image into Assets file as an “ARImage”



<https://developer.apple.com/videos/play/wwdc2018/610/>

object recognition in ARKit

- detection based on scans of an object
- scanning app available here:
 - https://developer.apple.com/documentation/arkit/scanning_and_detecting_3d_objects



Prepare to scan

Define bounding box

Scan

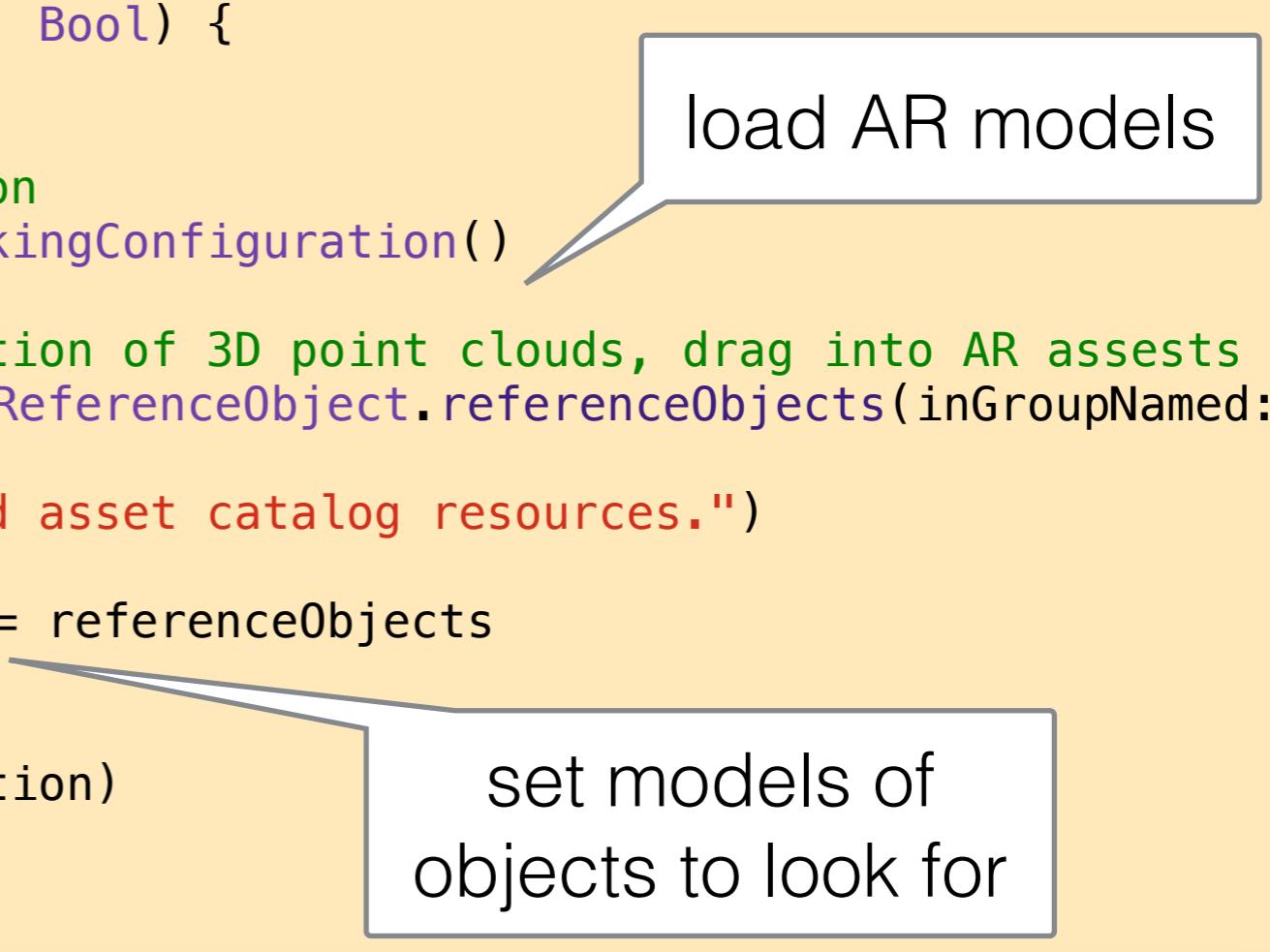
Adjust origin

Test and Export

object recognition in ARKit

- drag **aobject** into assets repository
- tell ARKit to begin looking for object

```
override func viewWillAppears(animated: Bool) {  
    super.viewWillAppear(animated)  
  
    // Create a session configuration  
    let configuration = ARWorldTrackingConfiguration()  
  
    // here is where we setup detection of 3D point clouds, drag into AR assets  
    guard let referenceObjects = ARReferenceObject.referenceObjects(inGroupNamed:  
        "gallery", bundle: nil) else {  
        fatalError("Missing expected asset catalog resources.")  
    }  
    configuration.detectionObjects = referenceObjects  
  
    // Run the view's session  
    sceneView.session.run(configuration)  
}
```



load AR models

set models of objects to look for

object recognition in ARKit

- use delegation to respond when object is found
- add new nodes anchored to the ARAnchor

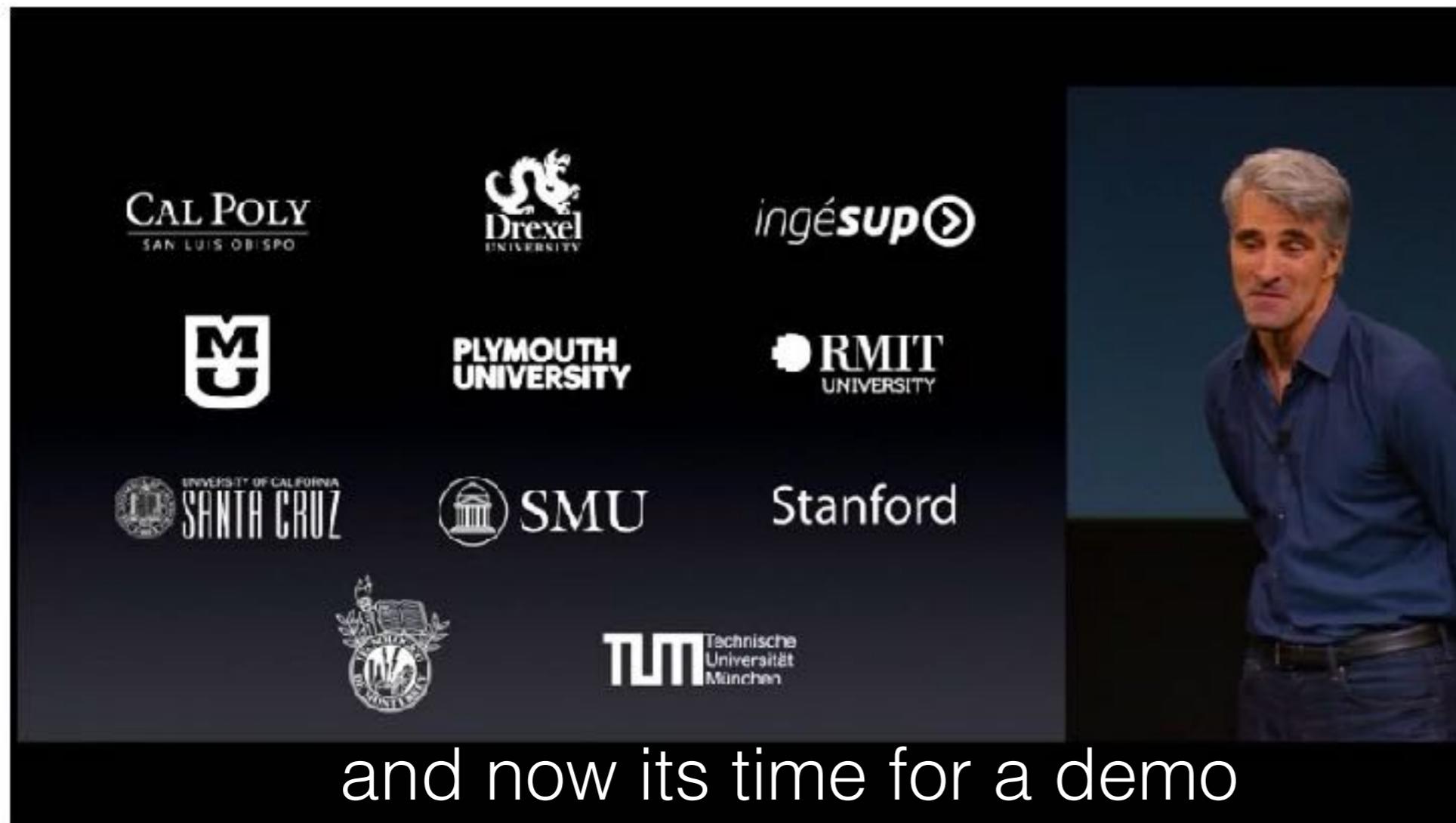
```
func renderer(_ renderer: SCNSceneRenderer,  
             didAdd node: SCNNode,  
             for anchor: ARAnchor) {  
  
    if let objectAnchor = anchor as? ARObjectAnchor {  
  
        newNode = ... what you want it to be ...  
        print(objectAnchor.name! + " found")  
        node.addChildNode(newNode)  
    }  
}
```

delegate function called
when new node is added to
ARSCNScene

true if recognized as
ARObject from gallery

ARKit with ObjectRecognition

- find the object!
- place images around the object!



A collage of university logos and a photograph of Steve Jobs. The logos include CAL POLY, Drexel UNIVERSITY, ingéSUP, PLYMOUTH UNIVERSITY, RMIT UNIVERSITY, UNIVERSITY OF CALIFORNIA SANTA CRUZ, SMU, Stanford, and TUM Technische Universität München. To the right of the logos is a photograph of Steve Jobs speaking on stage.

and now its time for a demo



Planes as Content Holders

- planes are placeholders for flat content like images and videos
- planes are the geometry for a node in SceneKit



```
// add in a video near the detected object
let planarNode = SCNNode()
let contentPlane = SCNPlane(width: CGFloat(0.1), height: CGFloat(0.05))
let avMaterial = SCNMaterial()
```

```
avMaterial.diffuse.contents = ...image or video...
contentPlane.materials = [avMaterial]
```

```
planarNode.geometry = contentPlane
node.addChildNode(planarNode)
```

setup image or video

set material of the plane

plane is geometry of node



Simple Nodes in SceneKit

- define geometry with pre-made shape, like box
- make new node with given geometry



```
let boxNode:SCNNode? = SCNNode()  
  
let box = SCNBox(width: CGFloat(0.1), height: CGFloat(0.1),  
                  length: CGFloat(0.1), chamferRadius: 0.01)  
box.firstMaterial?.diffuse.contents = UIColor(white: 1.0, alpha: 0.8)  
box.firstMaterial?.isDoubleSided = true  
  
boxNode!.geometry = box // make this node a box!  
  
node.addChildNode(box!)
```

setup geometry and color

```
let alphaAction = SCNAction.fadeOpacity(to: 0.1, duration: 5)  
box?.runAction(alphaAction)
```

add animation!



Actions in SceneKit

- any movement scaling, or change of added nodes should be done via actions
- actions can be reused on different nodes!

```
// add some animation  
let prevScale = uiObject.scale  
uiObject.scale = SCNVector3(CGFloat(prevScale.x)/3, CGFloat(prevScale.y)/3,  
                           CGFloat(prevScale.z)/3)
```

not yet on screen: will start smaller,
then scale back to original size

```
let scaleAction = SCNAction.scale(to: CGFloat(prevScale.x), duration: 2)  
scaleAction.timingMode = .easeIn
```

```
node.addChildNode(uiObject)
```

make visible

define scaling action and params

```
uiObject.runAction(scaleAction, forKey: "scaleAction")
```

run action immediately

UI Text Elements in SceneKit



```
func createTextNode(textString: String) -> SCNNode?{  
  
    let textNode: SCNNode? = SCNNode()  
    textNode!.geometry = setupTextParameters(textString: textString)  
    textNode!.scale = SCNVector3Make(0.001, 0.001, 0.001)  
    textNode!.position = SCNVector3Make(-0.1, 0.1, 0.0)  
    textNode?.castsShadow = true  
  
    return textNode  
}
```

lower left corner

```
func setupTextParameters(textString: String) -> SCNTex{  
    let text = SCNTex(string: textString, extrusionDepth: 1)  
  
    let material = SCNMaterial()  
    material.diffuse.contents = UIColor.white  
  
    text.flatness = 0  
    text.isWrapped = true  
    text.materials = [material]  
    return text  
}
```

if plane detection is enabled

extrusion defines depth of the text (flat versus extended)

color and other properties are set in the material of the node



<https://developer.apple.com/documentation/scenekit/scntext>

ARKit with ObjectRecognition

- add to demo some more interesting UI elements

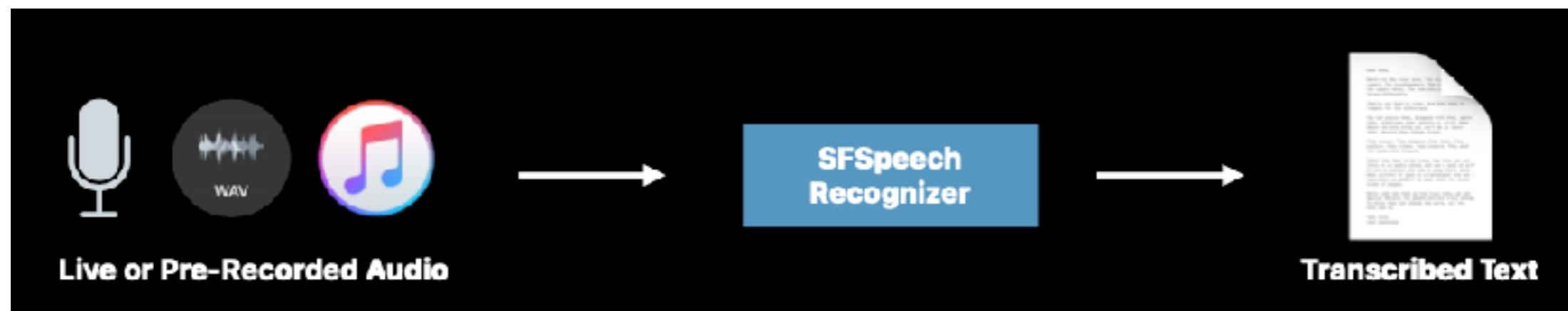


A dark grey rectangular area containing logos of nine universities: CAL POLY, Drexel UNIVERSITY, ingéSUP®, RMIT UNIVERSITY, UNIVERSITY OF CALIFORNIA SANTA CRUZ, SMU, Stanford, and TUM Technische Universität München. To the right of this area is a photograph of Steve Jobs speaking on stage.

and now its time for a demo

overview

- introduced in 2016, same technology underpinning Siri
 - user must provide explicit authorization
 - free, but limited to certain number of recognitions per day
 - only allowed to dictate about 1 minute of audio (siri)
 - supports streaming audio and file I/O



general usage

- uses API similar to REST (like Core Vision)
 - create a task
 - configure it (options)
 - start task
 - use completion handler (for updates and final text)
- best practices:
 - signify to user that the app is recording
 - show dictation as it happens

tradeoff server/on-device

- 2016: speech is translated via cloud services
- 2019: also available on device

	Server	On-device
Accuracy	Best	Good
Limits	1 minute max audio duration Limited requests per day	None
Languages	50+	10+

using SFSpeechRecognizer



- callback model
- needs AVFoundation for adding audio chunks

```
private let speechRecogniser = SFSpeechRecognizer(locale: Locale(identifier: "en-US"))!
private var recognitionRequest: SFSpeechAudioBufferRecognitionRequest?
private var recognitionTask: SFSpeechRecognitionTask?
private let audioEngine = AVAudioEngine()

let audioSession = AVAudioSession.sharedInstance()
audioSession.setCategory(AVAudioSession.Category.record)
audioSession.setMode(AVAudioSession.Mode.measurement)
audioSession.setActive(true, options: .notifyOthersOnDeactivation)

let inputNode = audioEngine.inputNode
let recordingFormat = inputNode.outputFormat(forBus: 0)

inputNode.installTap(onBus: 0, bufferSize: 1024, format: recordingFormat)
{ (buffer: AVAudioPCMBuffer, when: AVAudioTime) in
    self.recognitionRequest?.append(buffer)
}

audioEngine.prepare()
audioEngine.start()
```

much of the code also has **guards** for **error checking**

using SFSpeechRecognizer



```
let inputNode = audioEngine.inputNode
let recordingFormat = inputNode.outputFormat(forBus: 0)

inputNode.installTap(onBus: 0, bufferSize: 1024, format: recordingFormat)
{ (buffer: AVAudioPCMBuffer, when: AVAudioTime) in
    self.recognitionRequest?.append(buffer)
}

// perform on device, if possible
if speechRecogniser.supportsOnDeviceRecognition {
    recognitionRequest?.requiresOnDeviceRecognition = true
}

recognitionTask = speechRecogniser.recognitionTask(with: recognitionRequest)
{ [unowned self] result, error in
    if let result = result {
        let transcribedText = result.bestTranscription.formattedString
        // do something with text
    }
}

if result?.isFinal ?? (error != nil) {
    // this will remove the listening tap
    // so that the transcription stops
    inputNode.removeTap(onBus: 0)
}
```

more advanced speech processing

- recognition result contains many aspects of the voice, including:
 - Transcribed text (as we have seen)
 - Alternate transcriptions
 - Confidence in result
 - Timing
 - Speaking rate
 - Pause duration
 - Voice analytics

more advanced speech processing

```
bestTranscription=<SFTranscription>,
formattedString=I helped Apple recognize speech,
speakingRate=0.00000, averagePauseDuration=0.00000,
segments=(
    <SFTranscriptionSegment>, substringRange={0, 1}, timestamp=0.54, duration=0.24,
    confidence=0.966,
    substring=I, alternativeSubstrings=(\n),
    phoneSequence=AY,
    ipaPhoneSequence=\U02c8a\U0361\U026a, voiceAnalytics=null,
    <SFTranscriptionSegment>, substringRange={2, 6}, timestamp=0.78,
    duration=0.3600000000000001, confidence=0.966,
    substring=helped, alternativeSubstrings=(\n),
    phoneSequence=h EH l p t,
    ipaPhoneSequence=h.\U02c8\U025b.l.p.t,
    voiceAnalytics=null,
    ...
    <SFTranscriptionSegment>, substringRange={25, 31}, timestamp=2.49,
    duration=0.5699999999999998, confidence=0.966,
    substring=speech, alternativeSubstrings=(\n),
    phoneSequence=s p EE ch,
    ipaPhoneSequence=s.p.\U02c8i.t\U0361\U0283, voiceAnalytics=null
),
```

running on device will limit the available features!

more advanced speech processing

Each segment now has incredible amount of information

```
<SFTranscriptionSegment>,
substringRange={0, 10}, timestamp=0.27, duration=0.65, confidence=0.911,
substring=Performing,
alternativeSubstrings=(\n),
phoneSequence=(null),
ipaPhoneSequence=(null),

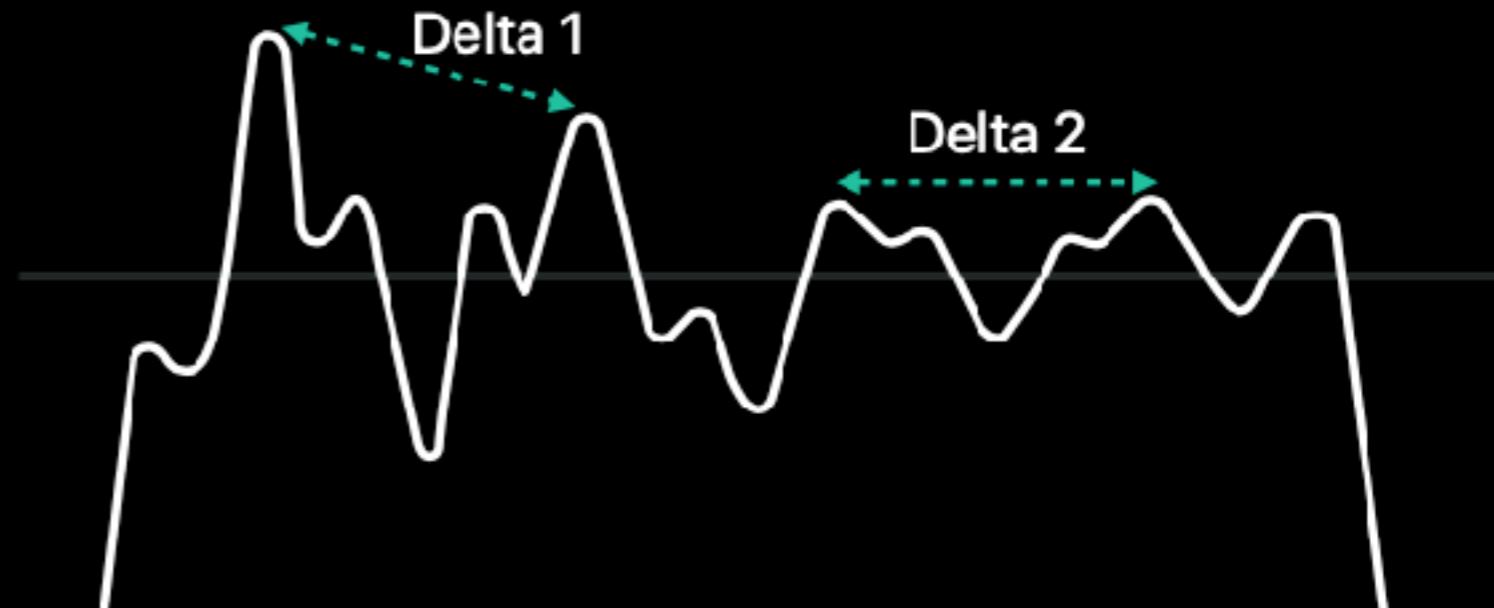
voiceAnalytics=<SFVoiceAnalytics>,
jitter=<SFAcousticFeature>, featureValues=(12.53122 ... 0.6218916),
frameDuration=0.010000,
shimmer=<SFAcousticFeature>, featureValues=(0.7158176 ... 2.518468),
frameDuration=0.010000,
pitch=<SFAcousticFeature>, featureValues=(0.8526305, ... 0.04258926),
frameDuration=0.010000,
voicing=<SFAcousticFeature>, featureValues=(0.07444749 ... 0.4056852),
frameDuration=0.010000",
```

each featureValues array is length of audio frames, could be **hundreds** of values

analytics

Jitter

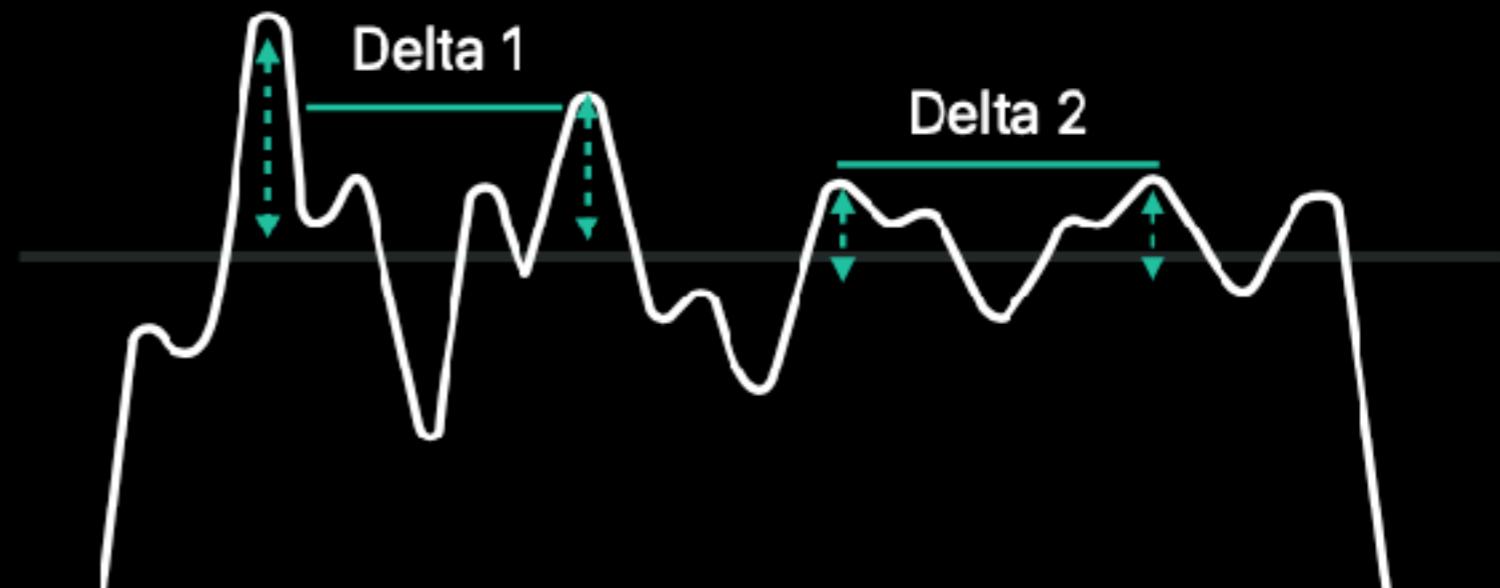
Measures variation in pitch



$$\text{Jitter} = \text{Delta1} - \text{Delta2}/\text{mean}$$

Shimmer

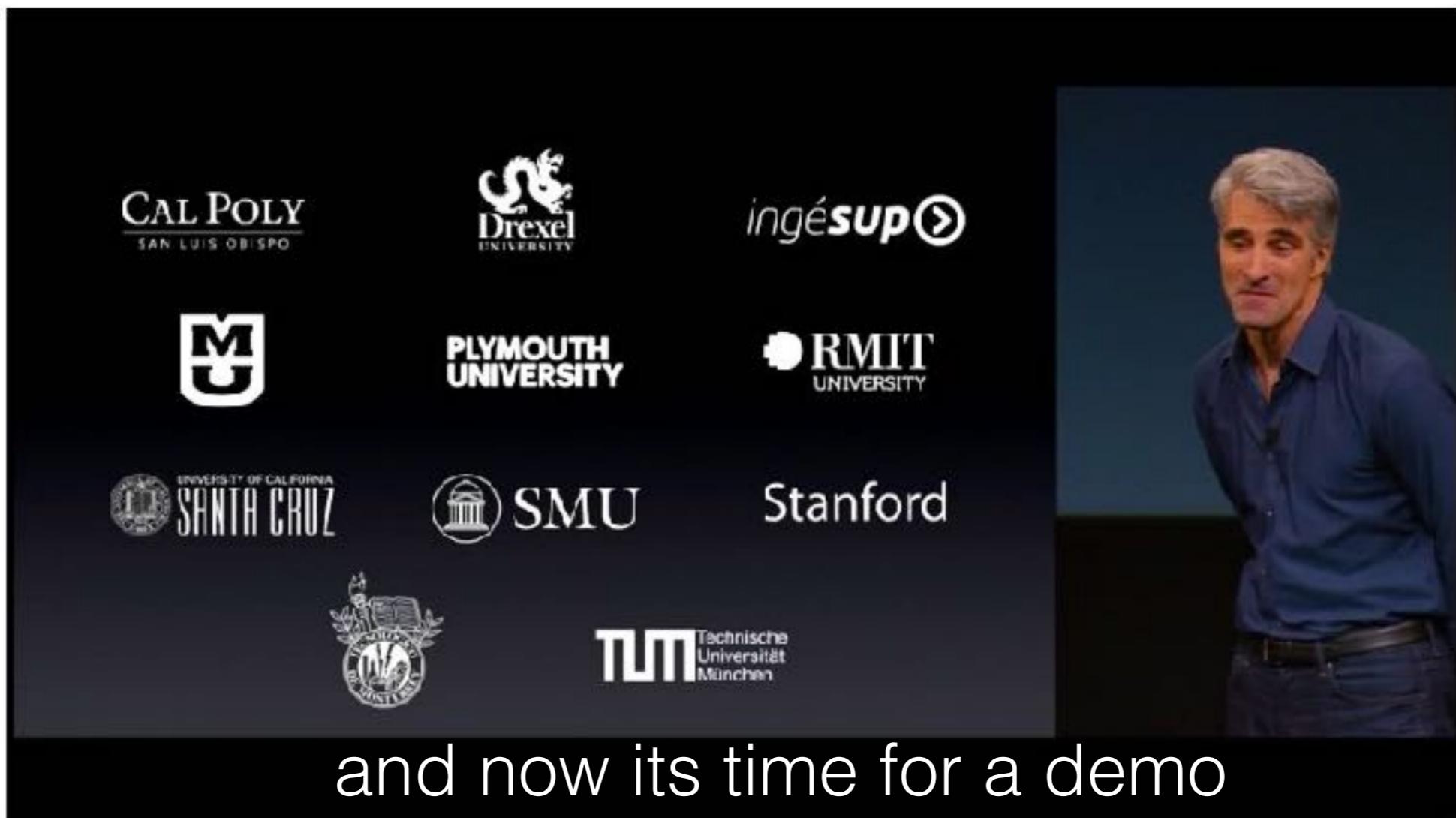
Measures variation in amplitude



$$\text{Shimmer} = \text{Delta1} - \text{Delta2}/\text{mean}$$

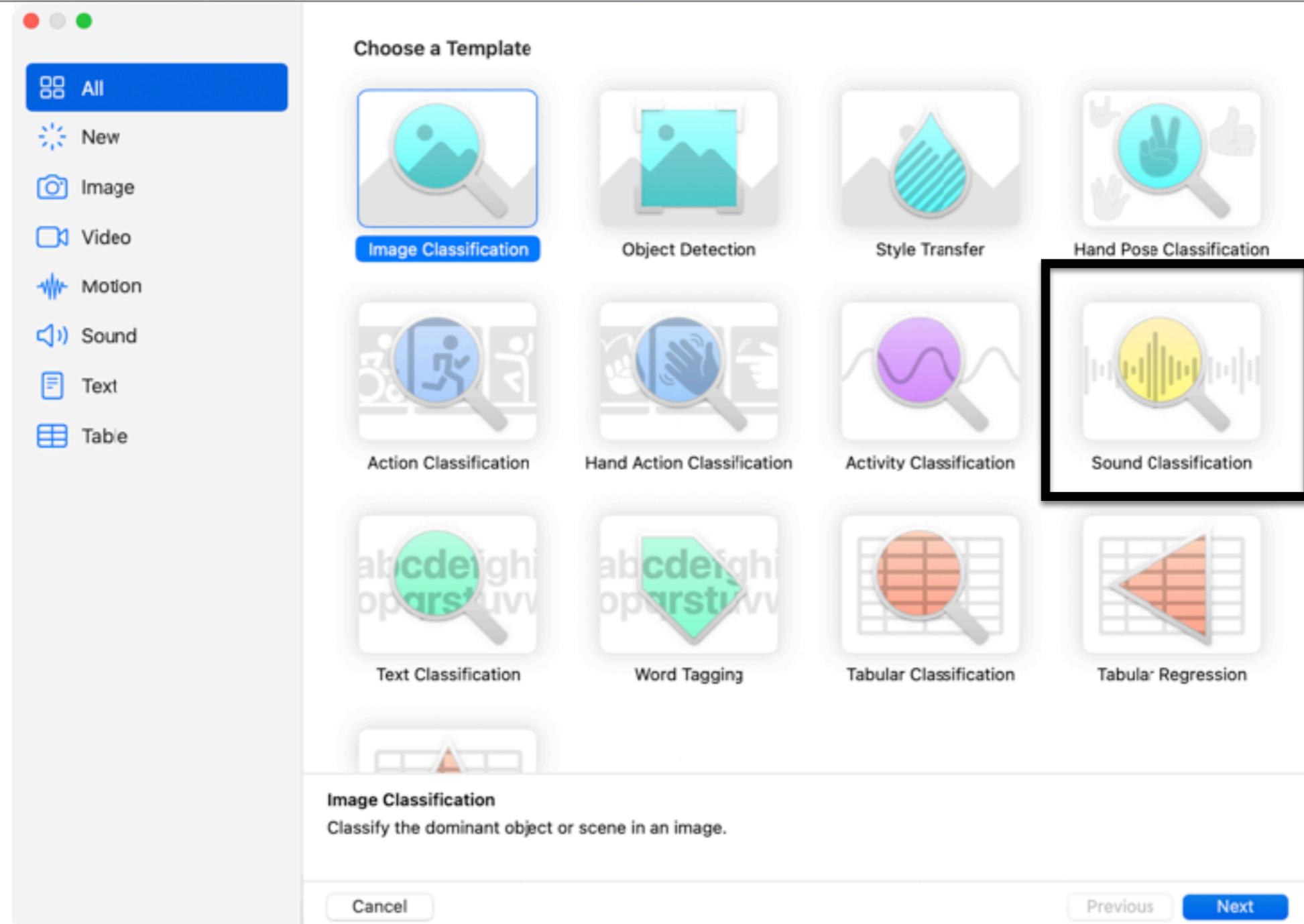
SFSpeechRecognizer

- adding audio blocks from input buffer

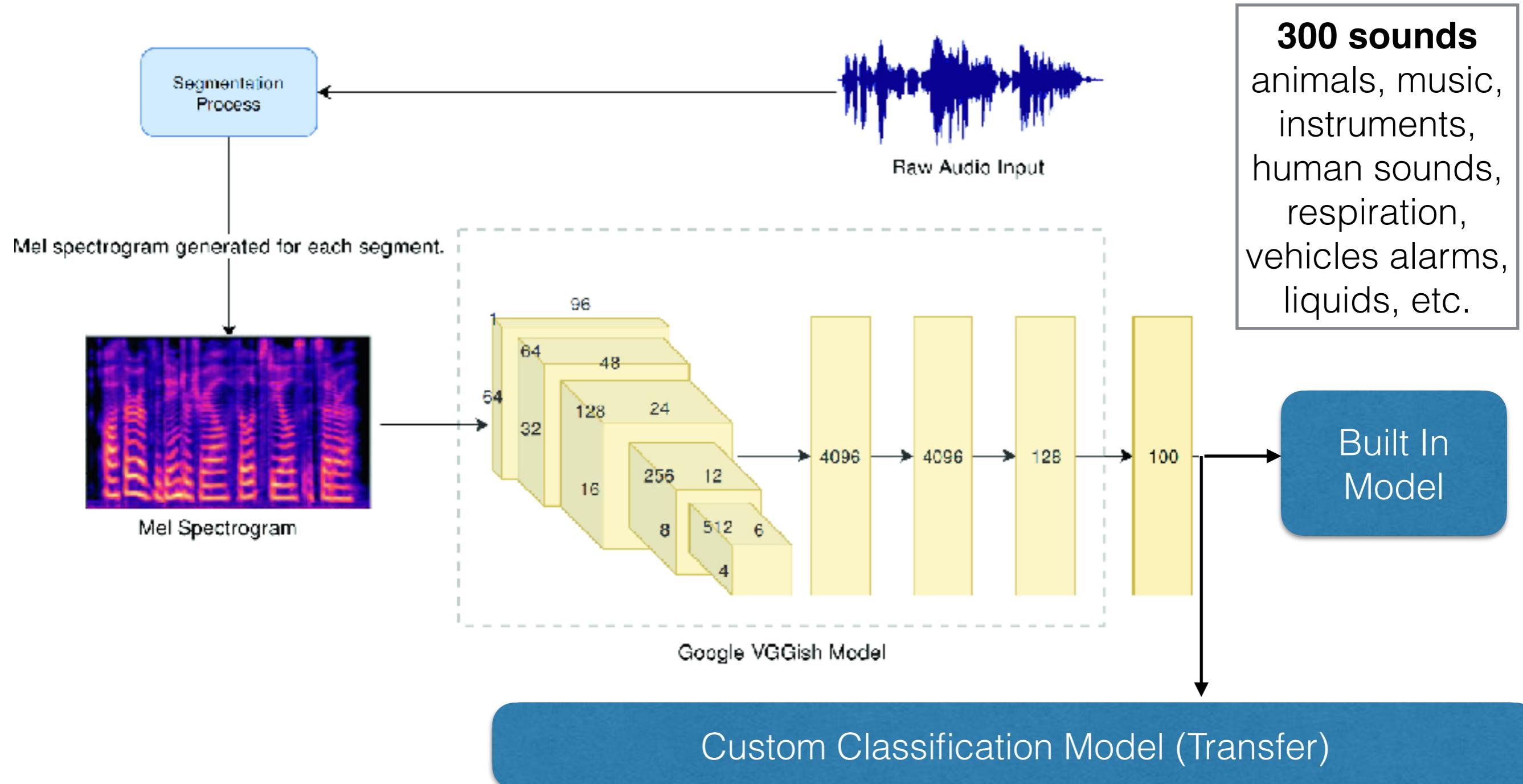


and now its time for a demo

Bonus: Create ML (iOS 15.0+)



Create ML Audio Analyzer



Create ML Audio Analyzer

The screenshot shows the 'MySoundClassifier.mlproj' project in the Create ML application. The project summary indicates 11 classes, 93% Training, 73% Validation, and 39% Testing, resulting in a 5.5 MB model file.

Data Inputs:

- Training Data:** 6,705 items from 'IRMAS-TrainingData'.
- Validation Data:** Auto.
- Testing Data:** 1,432 items from 'Part1'.

Parameters:

- Maximum Iterations: 25
- Overlap Factor: 50%

A message at the bottom states: "Training completed after 12 minutes, 16 seconds — today at 15:52". A "Make a Copy" button is also present.

<https://martinmitrevski.com/2019/12/09/sound-classification-with-create-ml-on-ios-13/>

Create ML Audio Analyzer

```
let request = try SNClassifySoundRequest(mlModel: soundClassifier.model)
try streamAnalyzer.add(request, withObserver: self)

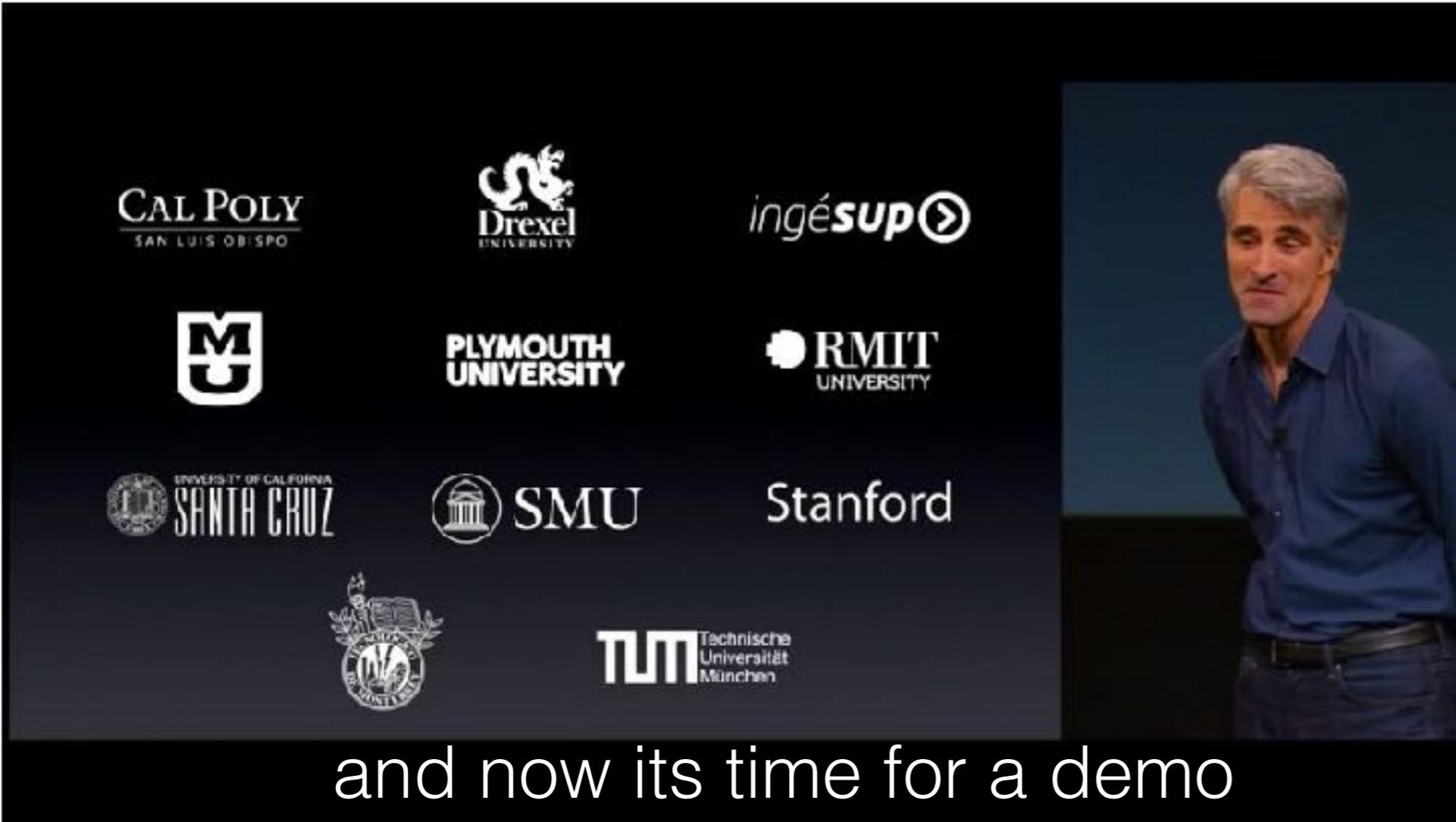
private func prepareForRecording() {
    let inputNode = audioEngine.inputNode
    let recordingFormat = inputNode.outputFormat(forBus: 0)
    streamAnalyzer = SNAudioStreamAnalyzer(format: recordingFormat)
    inputNode.installTap(onBus: 0, bufferSize: 1024, format: recordingFormat) {
        [unowned self] (buffer, when) in
        self.queue.async {
            self.streamAnalyzer.analyze(buffer, atAudioFramePosition: when.sampleTime)
        }
    }
    audioEngine.prepare()
    do { try audioEngine.start() } catch {...}\
}

func request(_ request: SNRequest, didProduce result: SNResult) {
    guard let result = result as? SNClassificationResult else { return }
    for classification in result.classifications {
        print(classification.identifier, classification.confidence)
    }
}
```

<https://martinmitrevski.com/2019/12/09/sound-classification-with-create-ml-on-ios-13/>

Sound Analysis

- add sound classification to our project

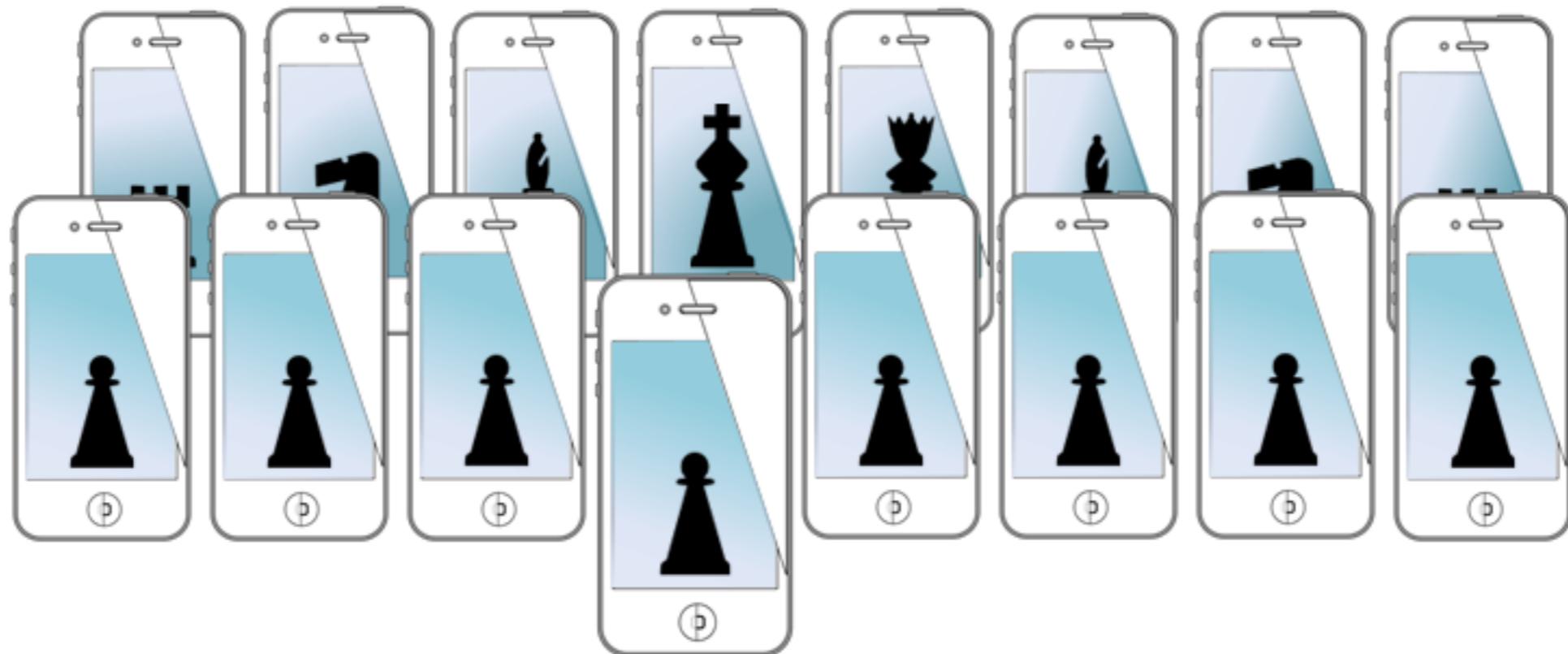


clock	⌚
bowling_impact	🎳
elk_bugle	🐑
bird_vocalization	🐦🎶
percussion	🥁
door_slam	🚪
foghorn	嘟嘟
horse_neigh	🐴��
applause	👏
writing	✍
train_horn	🚂嘟
fireworks	🎆
person_running	🏃
ukulele	🎸🌺
playing_hockey	🏒
mosquito_buzz	🦟
whoosh_swoosh_swish	💥💨
organ	🎹
typewriter	⌨️
alarm_clock	⏰
chopping_wood	⽊斧
bird_chirp_tweet	🐦🎶
theremin	🎶
ambulance_siren	救护车
water_pump	水泵
cow_moo	🐮哞
engine_idling	🚗喇叭

for next time...

- Pitching
- ~Fin~

MOBILE SENSING LEARNING



CS5323 & 7323
Mobile Sensing and Learning

Adding Vision Object Detection

Eric C. Larson, Lyle School of Engineering,
Computer Science, Southern Methodist University

Back Up Slides

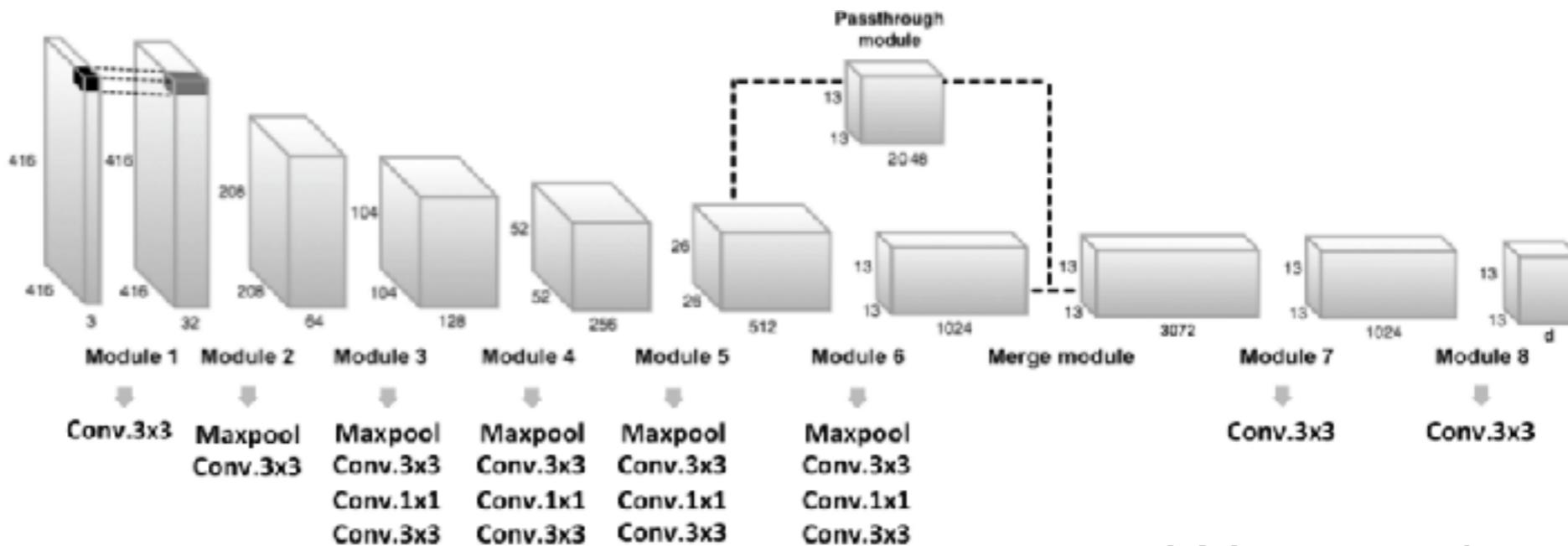
Detecting Objects in Images

- Using Turi Create
 - SFrames
 - Simple Machine Learning
- Object Detection with YOLO
 - Creation
 - Exporting to CoreML
 - Incorporating into ARKit

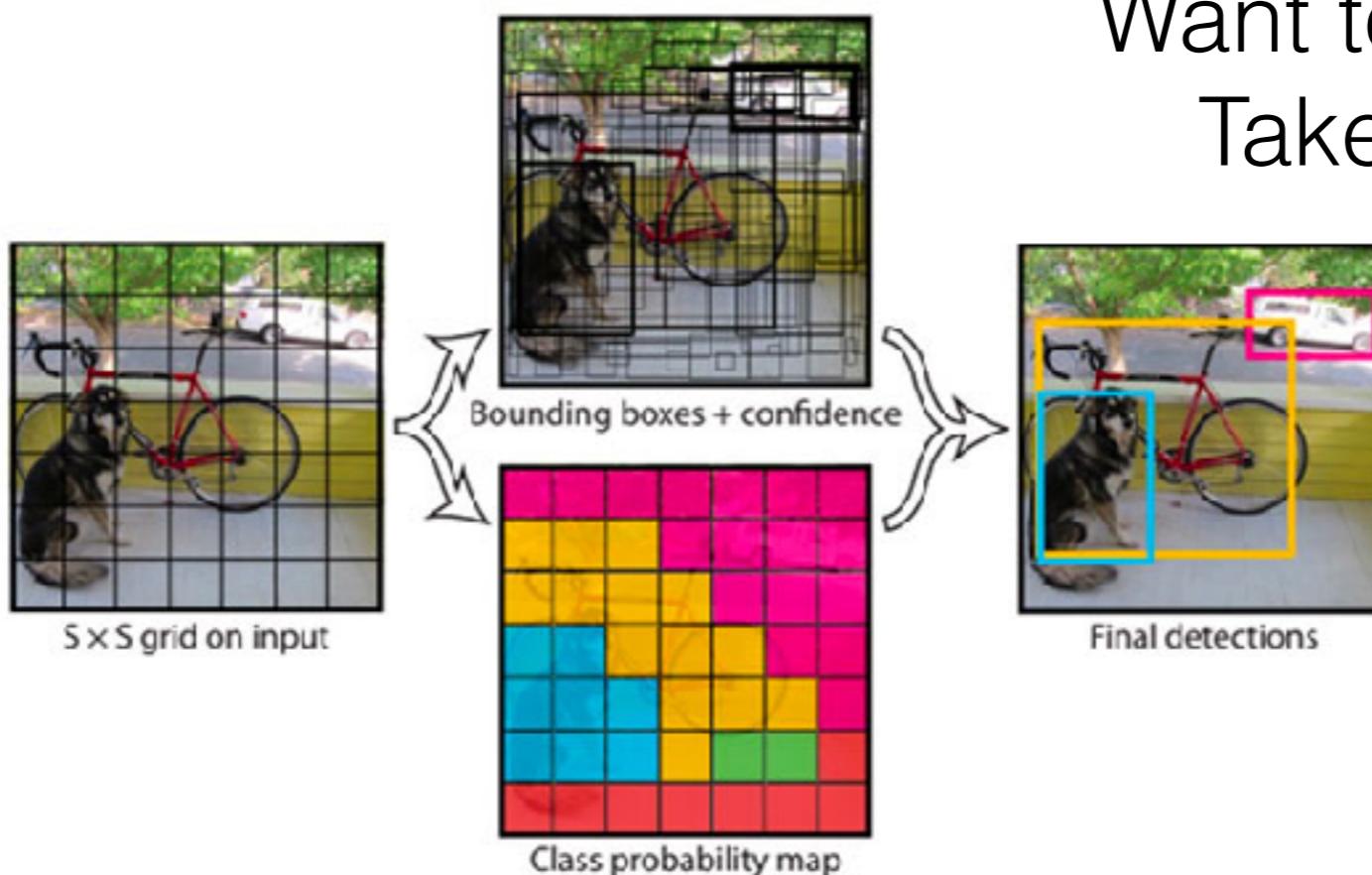
didn't we already detect objects?

- with ARKit, yes
- but what if we want to detect all sorts of objects
- and we need to do it very fast (multiple times per second)
- and we only have examples of the types of objects we want, not a scan of the exact object
- enter... Turi Create

YOLO



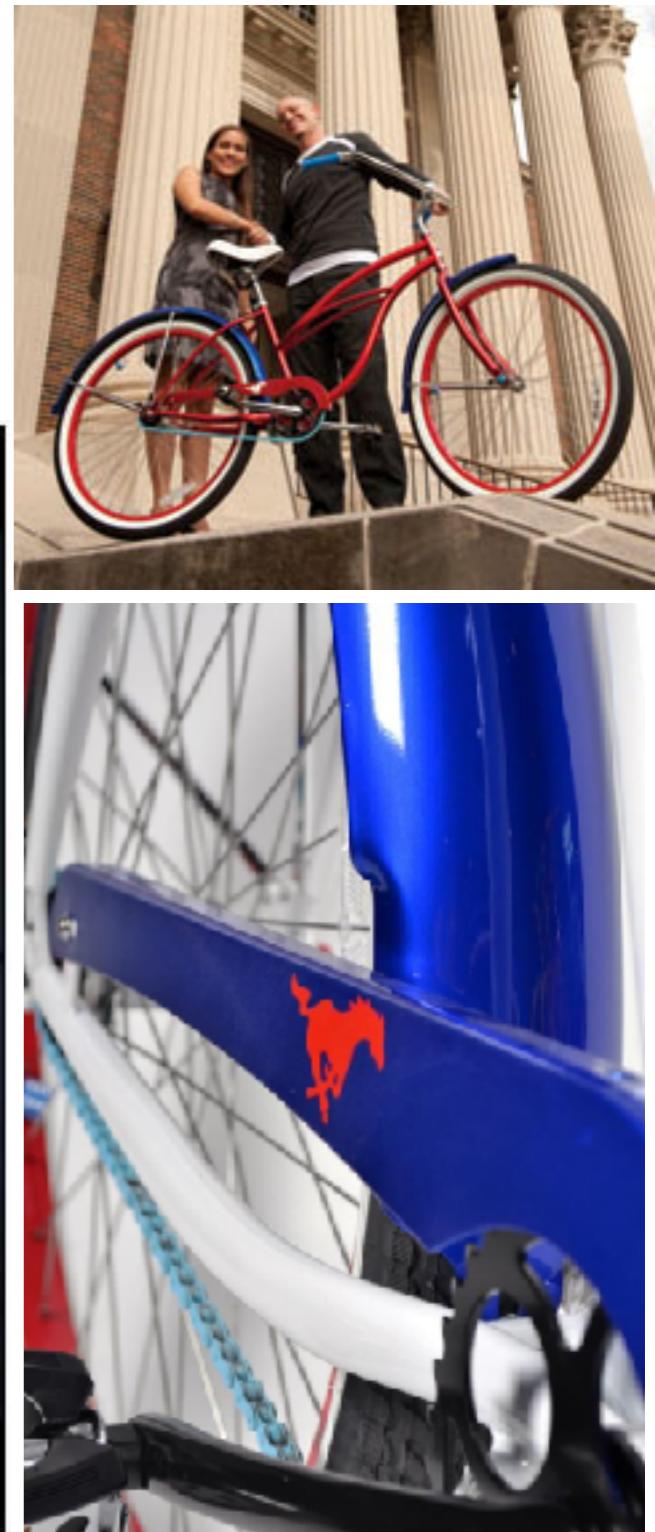
Want to know more?
Take CS8321...



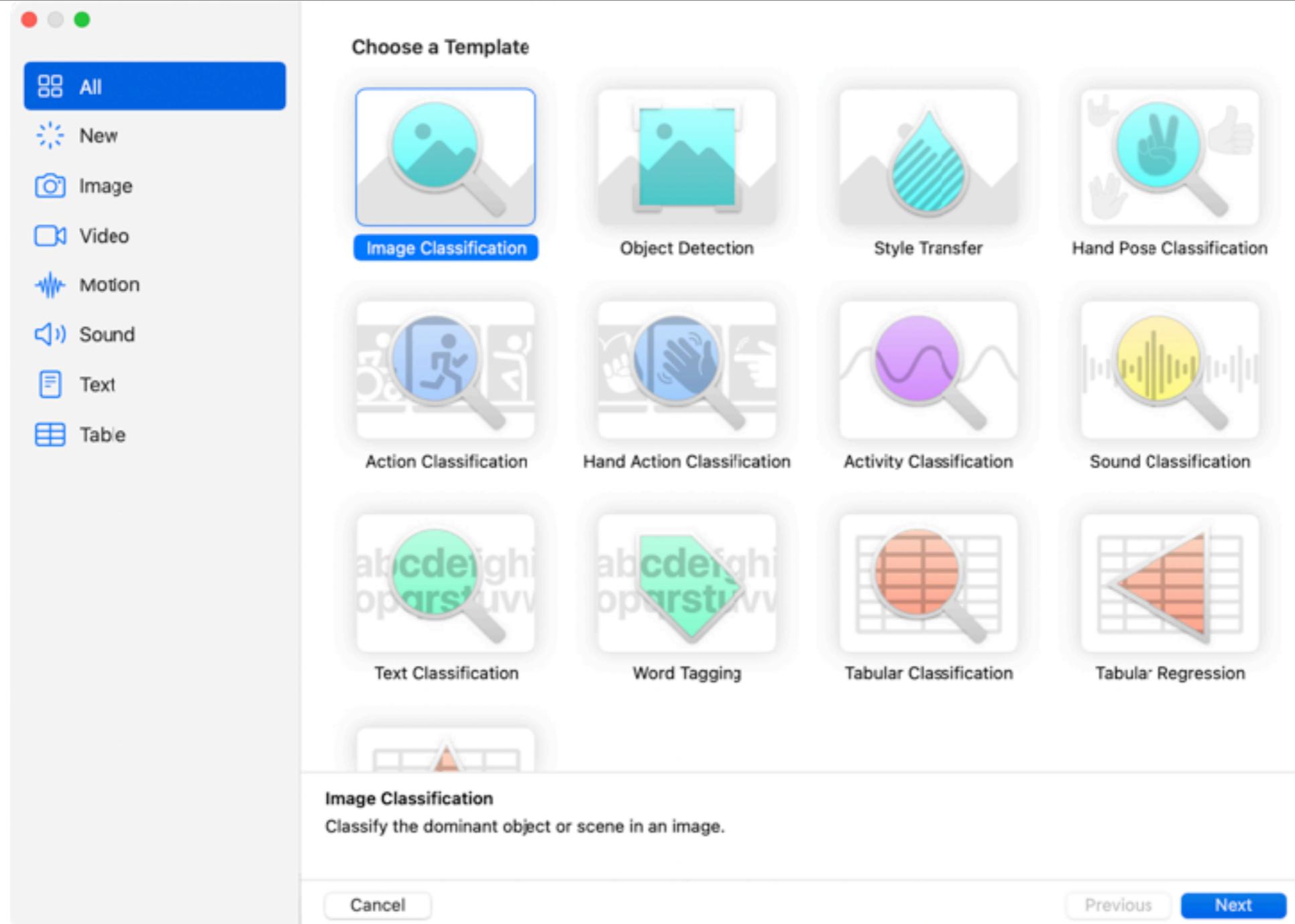
<https://datascience.stackexchange.com/questions/42509/yolo-algorithm-understanding-training-data>

create an object detector with Turi

- Turi Create, high level overview
 - SFrames and Object Detectors



Bonus: Create ML (iOS 15.0+)



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

Create ML Stylizer

The screenshot shows a user interface for creating an ML Stylizer. The top navigation bar includes 'Train More', 'Snapshot', 'Settings', **Training**, 'Preview', 'Output', and 'Activity'. The 'Training' tab is active.

Stylized Validation at Iteration 700: A stylized image of a dog's face, where the content image has been transformed to match the style of the 'Style' image.

Style: An image of Katsushika Hokusai's 'The Great Wave off Kanagawa'.

Loss vs Iterations: A line graph showing the Style Loss decreasing from approximately 22.510 at iteration 0 to about 12.915 at iteration 700. The Y-axis ranges from 0 to 1000, and the X-axis ranges from 0 to 200.

Content Loss vs Iterations: A line graph showing the Content Loss remaining relatively stable around 12.915 across 700 iterations. The Y-axis ranges from 0 to 20, and the X-axis ranges from 0 to 200.

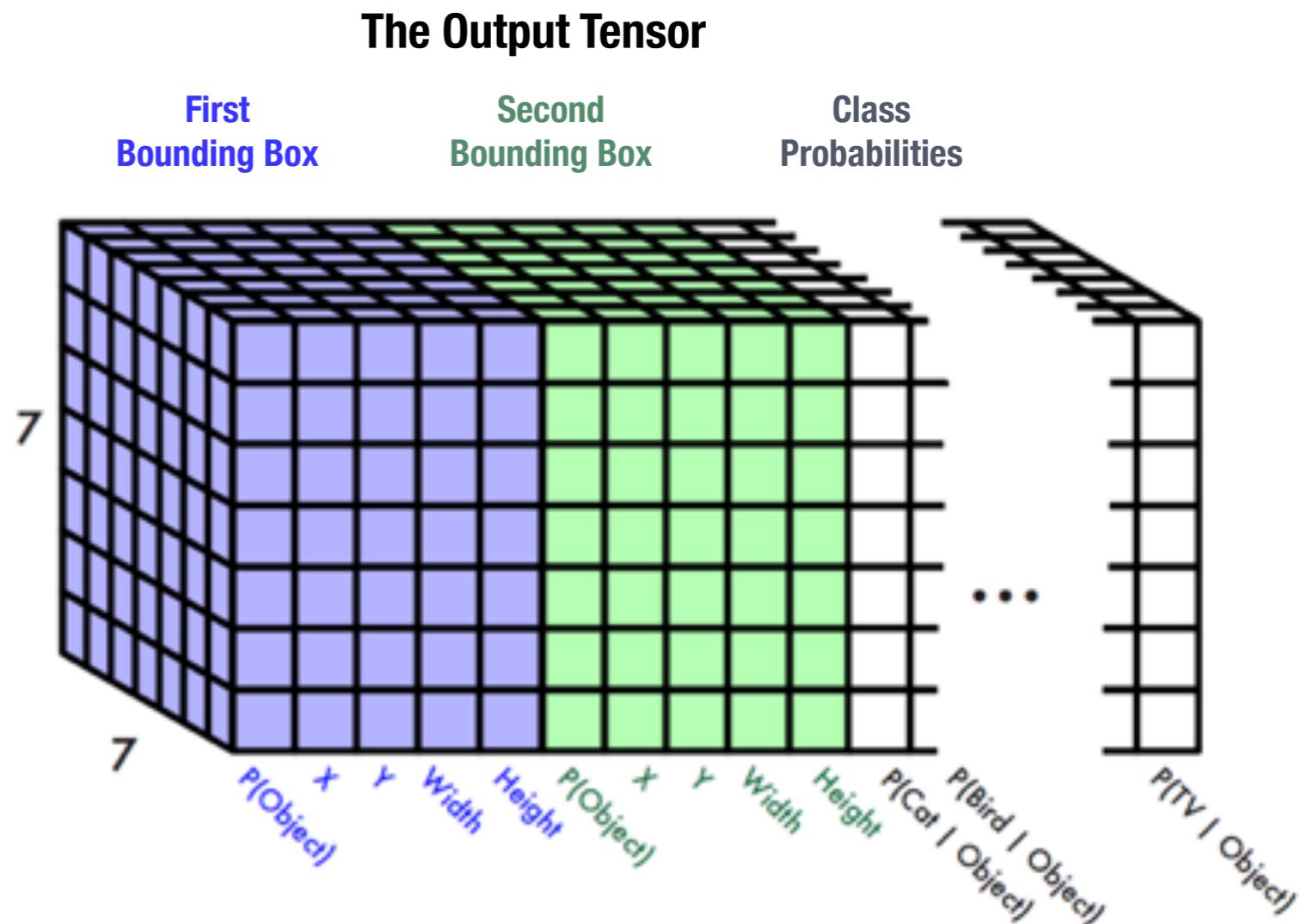
Status: A message at the bottom left indicates 'Completed 700 iterations'.

Activity Log: A list of events on the right side:

- Activity: Training Completed, Date: Nov 10, 2021, Time: 6:35 PM
- Activity: 700 Iterations
- Activity: Training Extended, Date: Nov 10, 2021, Time: 6:26 PM
- Activity: 200 more iterations
- Activity: Snapshot Saved, Date: Nov 10, 2021, Time: 6:26 PM
- Activity: Snapshot created at Iteration 500
- Activity: Training Started, Date: Nov 10, 2021, Time: 4:54 PM
- Activity: 500 Iterations
- Activity: Validation Data Added, Date: Nov 10, 2021, Time: 4:54 PM
- Activity: doge.jpg
- Activity: Training Data Added, Date: Nov 10, 2021, Time: 4:52 PM
- Activity: wave_style.png
- Activity: Content Images Added, Date: Nov 10, 2021, Time: 4:51 PM
- Activity: NaturalContentDataset
- Activity: Model Source Created, Date: Nov 10, 2021, Time: 4:50 PM
- Activity: ExampleStyleTransfer 1
- Activity: Project Created, Date: Nov 10, 2021, Time: 4:50 PM
- Activity: ExampleStyleTransfer

YOLO output tensor

- is there an object?
- where is the object?
- how large is the Object?
- what is the object?
- if competing, which
- object is most likely?



using the vision API

- Overview:
 - streaming video will require constant images from the AR session, **use delegation**
 - inside delegate, need to **grab** ARFrame, **convert** to pixel buffer, and **process** on a **background** thread
 - **handle vision** results in a new function
 - **update** the UI on the **main thread**



setting up vision

```
let model = PersonBike()  
private var requests = [VNRequest]()  
  
@discardableResult  
func setupVision(useCPUOnly:Bool) {  
  
    //MARK: One, Setup Vision  
    let visionModel = try VNCoreMLModel(for: model.model)  
  
    // use this request to setup the object recognition with Vision  
    let objectRecognition = VNCoreMLRequest(model: visionModel,  
                                              completionHandler:  
                                                self.handleObjectRecognitionResult)  
  
    objectRecognition.imageCropAndScaleOption = .scaleFill  
    objectRecognition.usesCPUOnly = [True / False]  
    self.requests = [objectRecognition]  
}  
  
load model from Turi  
setup vision wrapper  
save this request for later
```



using vision with AR

ARDelegate function,
called at 60 FPS

```
func renderer(_ renderer: SCNSceneRenderer, updateAtTime time: TimeInterval) {  
    //MARK: Two, Get Image Frame from AR  
    guard let frame = sceneView.session.currentFrame else { return }  
  
    guard self.currentBuffer == nil, else { return } // limit FPS of detection  
    self.currentBuffer = frame.capturedImage // the pixels to process  
  
    // run in the background so that AR doesn't suffer performance  
    DispatchQueue.global(qos: .background).async {  
        // setup input image for the request  
        let imageRequestHandler = VNImageRequestHandler(  
            cvPixelBuffer: self.currentBuffer,  
            orientation: ORIENTATION,  
            options: [:])  
  
        imageRequestHandler.perform(self.requests)  
    }  
}
```

start processing request
that we setup previously

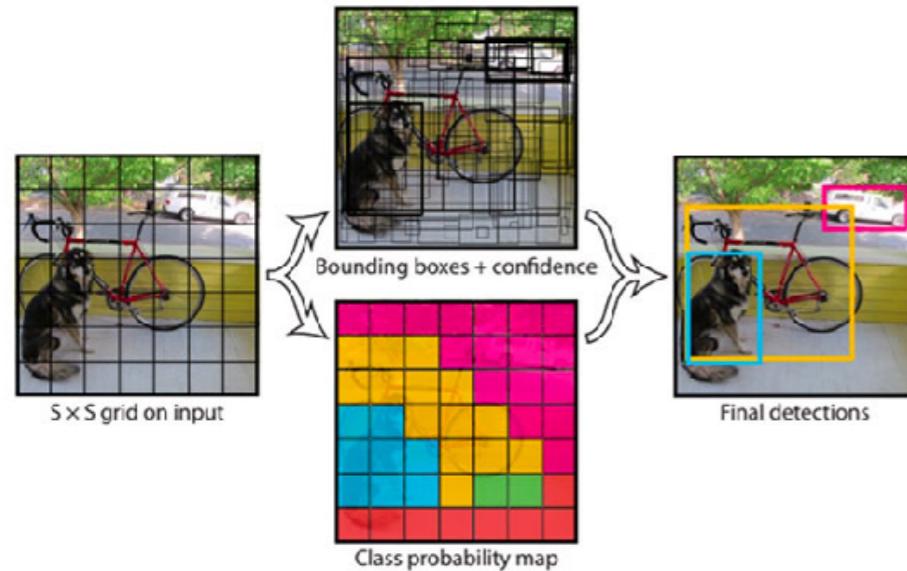
limit FPS of detection

setup image buffer

tell request what pixels to
process



handling output request



YOLO Output:

- list of bounding boxes for every square in grid
- list of all possible classes (bike, person)

can have multiple objects in the scene

```
func handleObjectRecognitionResult(_ request: VNRequest, error: Error?) {
    // perform all the UI updates on the main queue
    if let results = request.results { // if we have valid results, else its nil
        DispatchQueue.main.async(execute: {
            self.drawVisionRequestResults(results)
            self.updateOverlay() —————— go over, if time!
            // set as nil so we can process next ARFrame Image
            self.currentBuffer = nil —————— process the next
        })                                captured frame now
    }
}
```

object detection

- YOLO demo
- If time, go over creating overlays with CATransactions



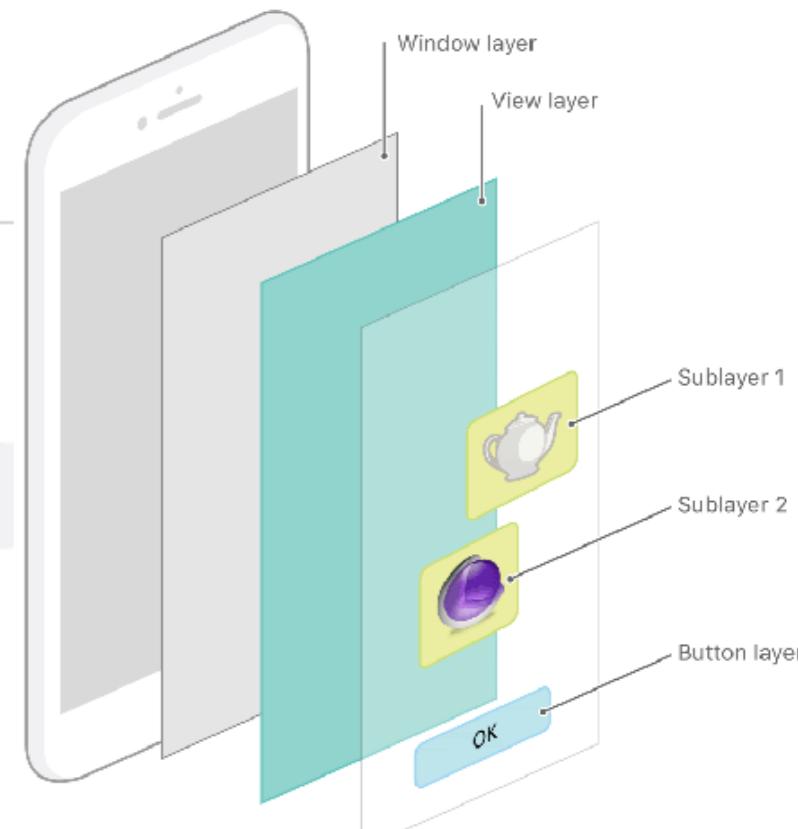
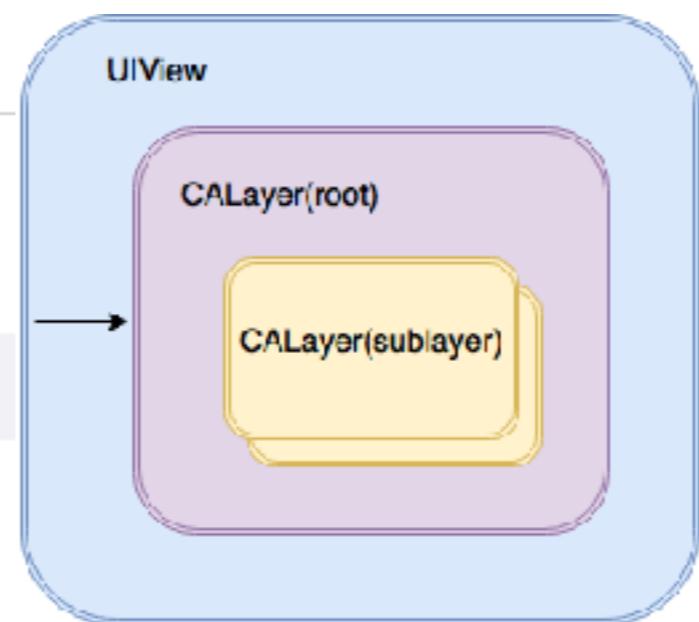
and now its time for a demo

CATransactions

A mechanism for grouping multiple layer-tree operations into atomic updates to the render tree.

Declaration

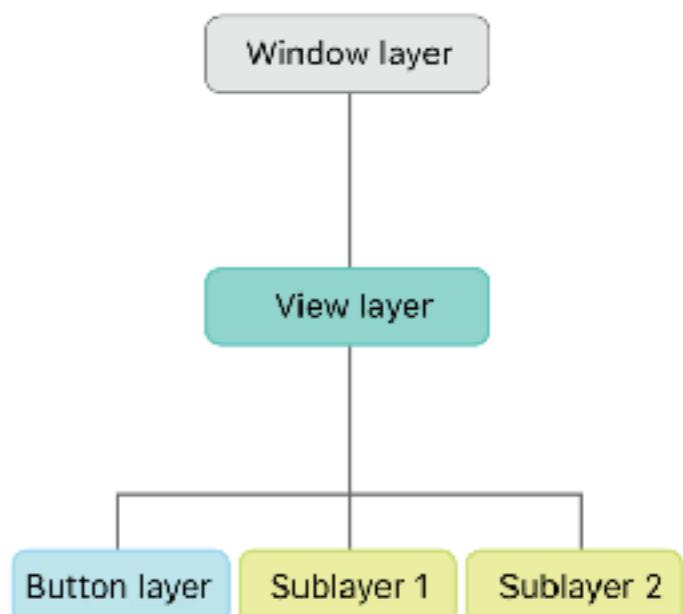
```
class CATransaction : NSObject
```



Overview

CATransaction is the Core Animation mechanism for batching multiple layer-tree operations into atomic updates to the render tree. Every modification to a layer tree must be part of a transaction. Nested transactions are supported.

Core Animation supports two types of transactions: *implicit* transactions and *explicit* transactions. Implicit transactions are created automatically when the layer tree is modified by a thread without an active transaction and are committed automatically when the thread's runloop next iterates. Explicit transactions occur when the application sends the **CATransaction** class a **begin()** message before modifying the layer tree, and a **commit()** message afterwards.



Overlays from YOLO

```
detectionOverlay = CALayer() // container layer that has all the renderings
detectionOverlay.name = "DetectionOverlay"

// set the initial bounds, will transform when we know more about the image
detectionOverlay.bounds = CGRect(x: 0.0, y: 0.0,
                                 width: self.view.bounds.width,
                                 height: self.view.bounds.height )

self.sceneView.layer.addSublayer(detectionOverlay)

func updateOverlay() {

    let bounds = self.view.bounds

    // ... Some magic transforms to the view ...
    // this tries to get the best mapping we can from the cropping that
    // Core Vision used, but It may not be 100% perfect
    // Scales, and adds magic numbers to the X and Y positions

    detectionOverlay.setNeedsDisplay() // sets display for all subviews
}
```

Overlays from YOLO

