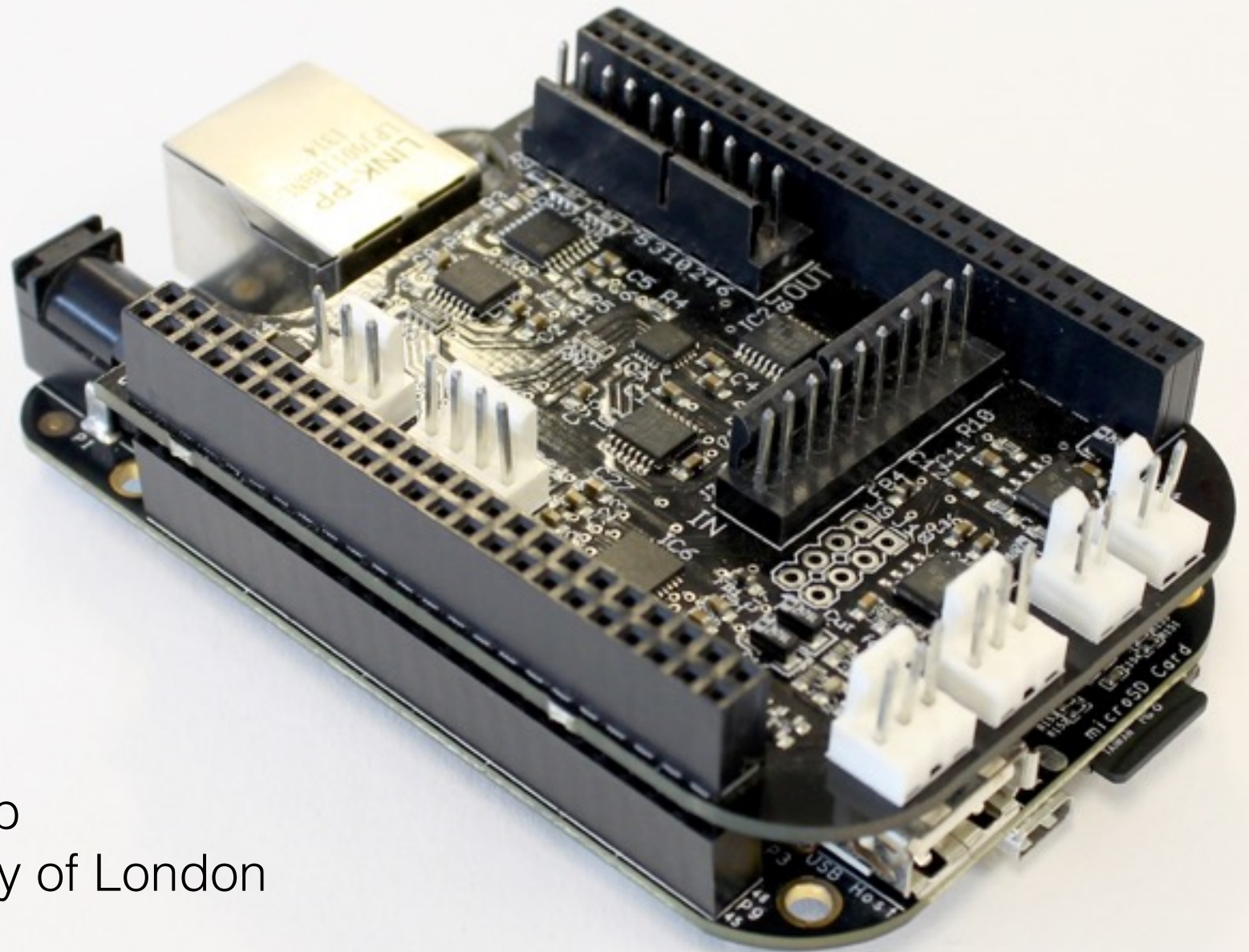


bela

*Ultra-low latency audio and
sensor processing
on the BeagleBone Black*



A project by
The Augmented Instruments Lab
at C4DM, Queen Mary University of London

<http://bela.io>

The Goal:

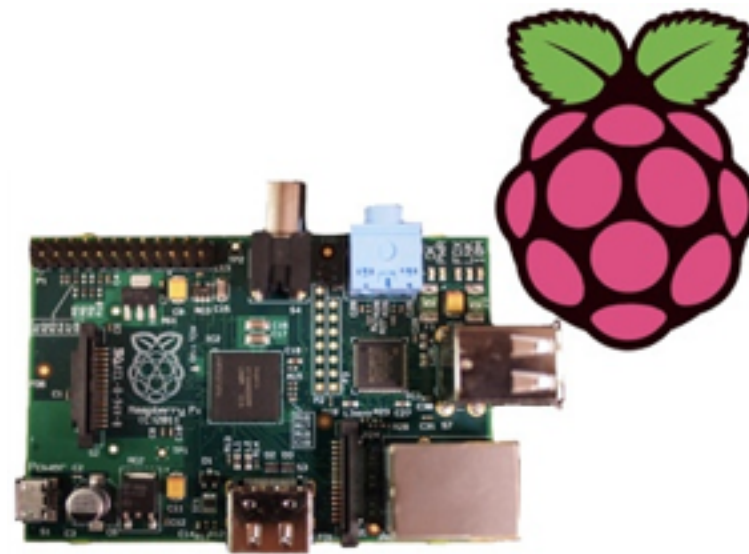
High-performance, self-contained
audio and sensor processing

The Goal:

High-performance, self-contained audio and sensor processing



- Easy low-level hardware connectivity
- No OS = precise control of timing
- Very limited CPU (8-bit, 16MHz)
- Not good for audio processing



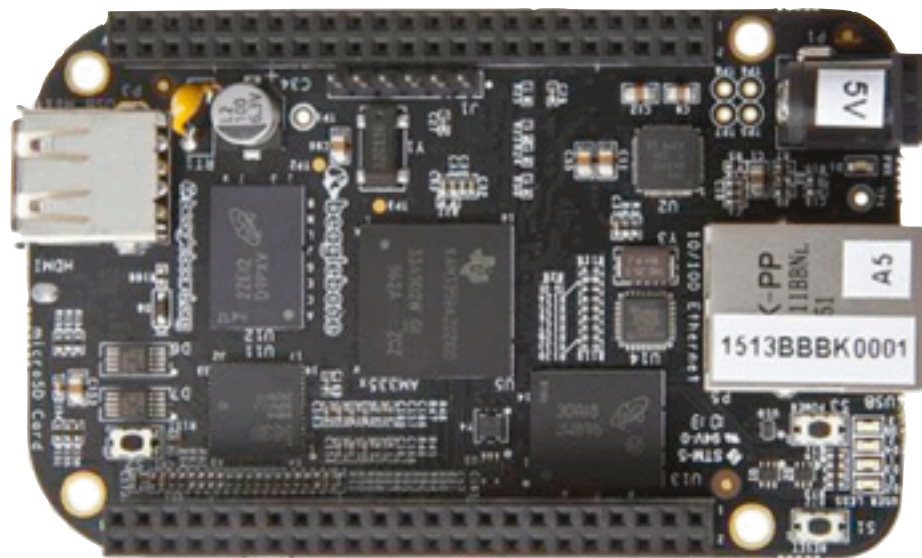
- Reasonable CPU (up to 1GHz ARM)
- High-level hardware (USB, network etc.)
- Limited low-level hardware
- Linux OS = high-latency / underruns



- Fast CPU
- High-level hardware (USB, network etc.)
- Arduino for low-level
- USB connection = high-latency, jitter
- Bulky, not self-contained

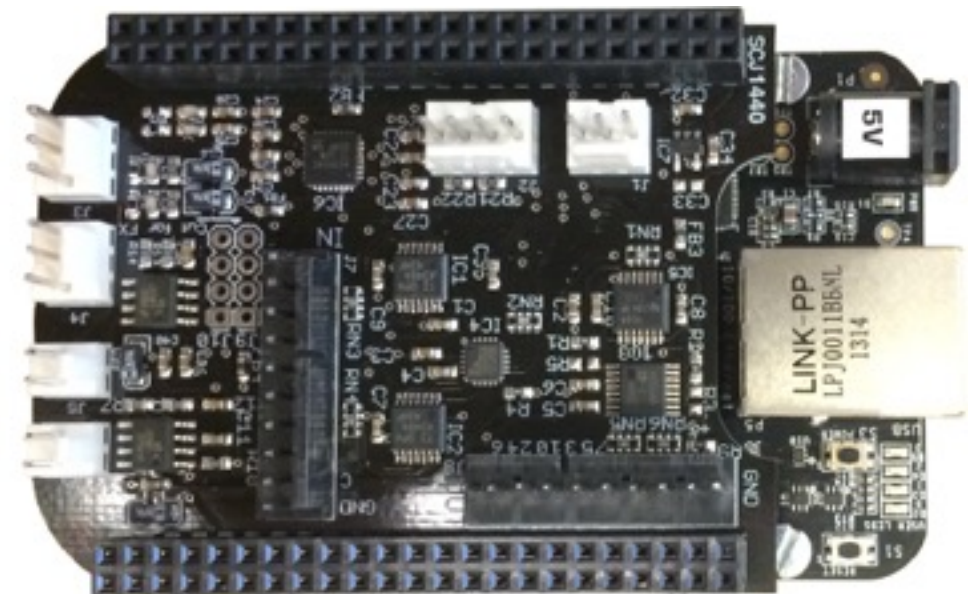
bela

hardware



BeagleBone Black

1GHz ARM Cortex-A8
NEON vector floating point
PRU real-time microcontrollers
512MB RAM



Custom Bela Cape

Stereo audio in + out
Stereo 1.1W speaker amps
8x 16-bit analog in + out
16x digital in/out



features

1ms round-trip audio latency without underruns

High sensor bandwidth: digital I/Os sampled at 44.1kHz; analog I/Os sampled at 22.05kHz

Jitter-free alignment between audio and sensors

Hard real-time audio+sensor performance, but full Linux APIs still available

Programmable using **C/C++, Pd or Faust**

Designed for **musical instruments and live audio**

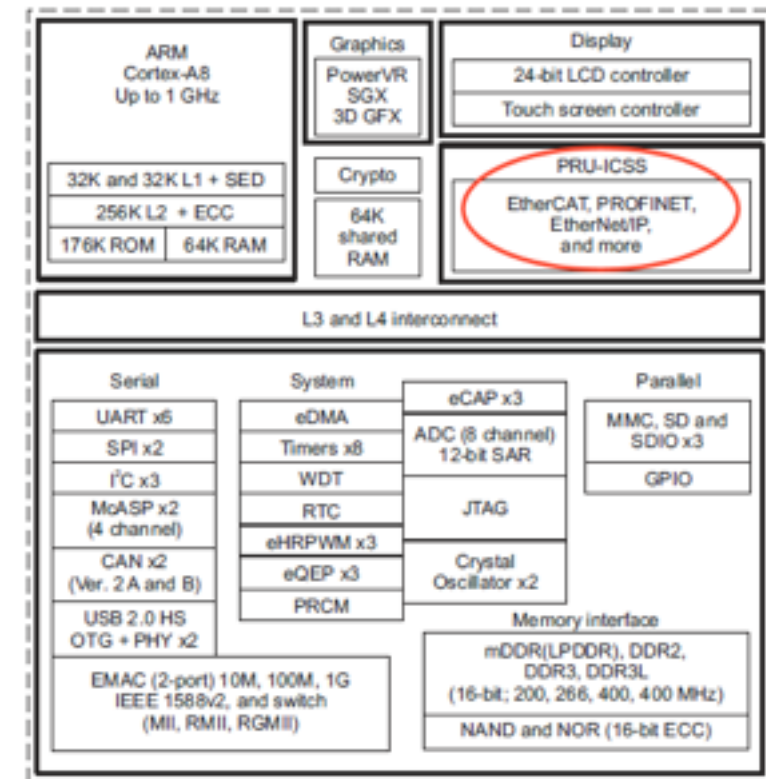
beld

software



Xenomai Linux kernel

Debian distribution
Xenomai hard real-time
extensions



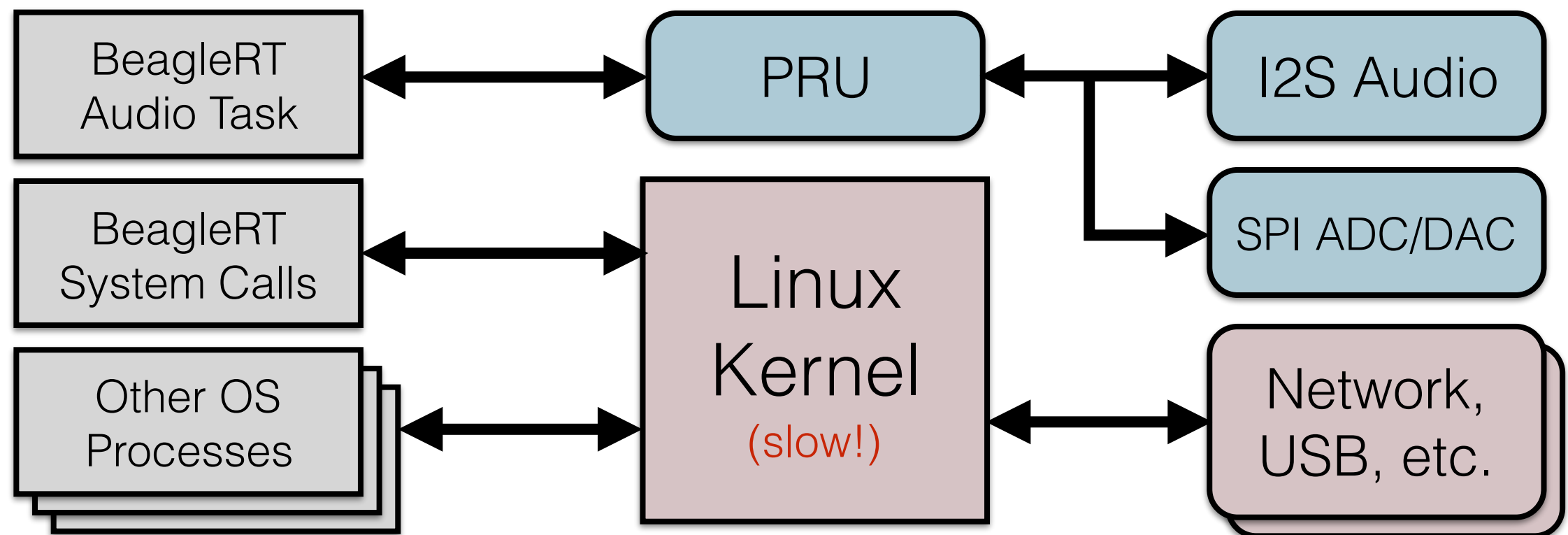
C++ programming API

Uses PRU for audio/sensors
Runs at higher priority
than kernel = *no dropouts*
Buffer sizes as small as **2**

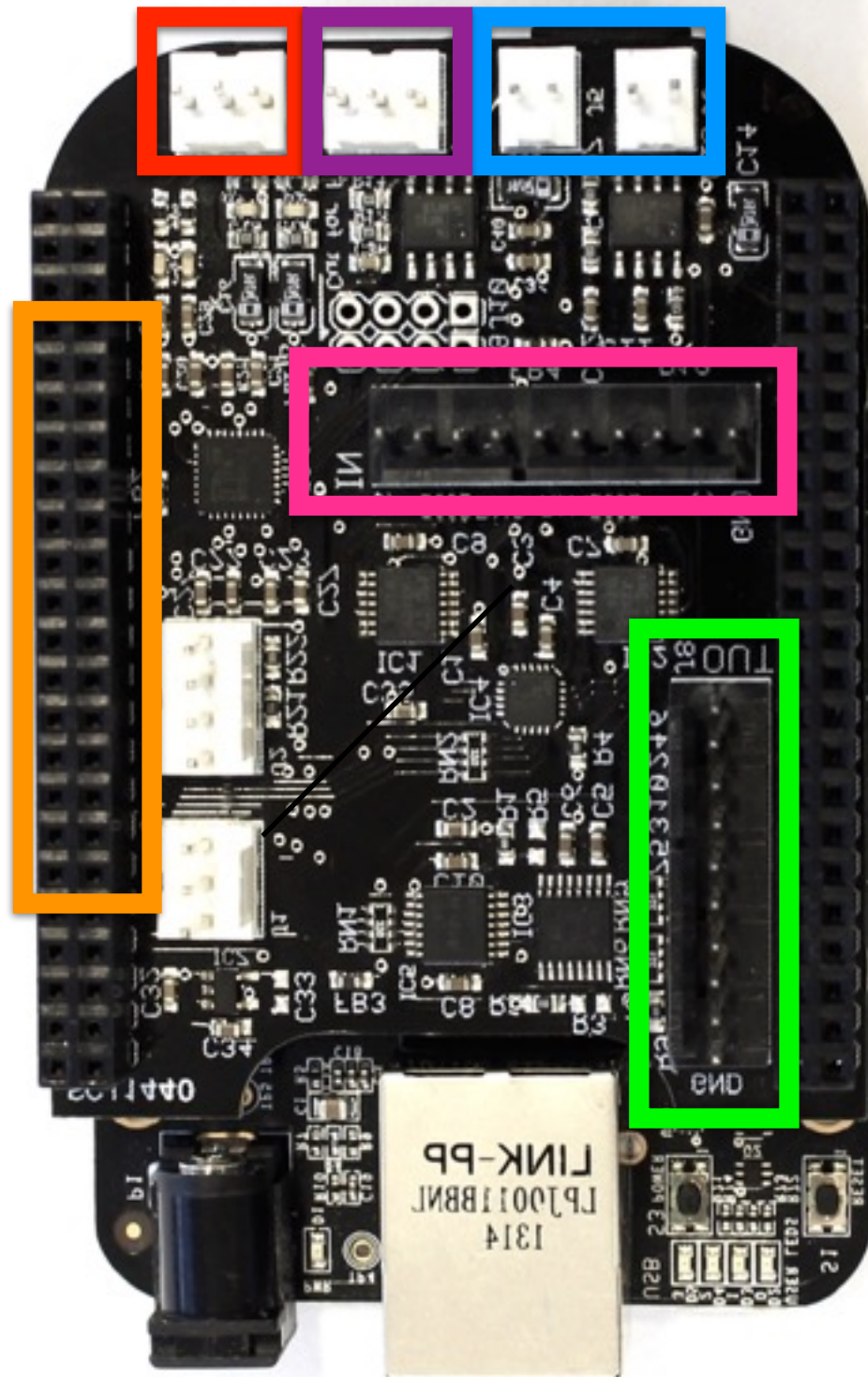
Bela software



- **Hard real-time** environment using Xenomai Linux kernel extensions
- Use BeagleBone **Programmable Realtime Unit (PRU)** to write straight to hardware



- Sample all matrix ADCs and DACs at **half audio rate** (22.05kHz)
- Buffer sizes as small as **2 samples** (90µs latency)



- **Speakers** with on-board amps
- **Audio In**
- **Audio Out**
- **16x digital I/O**
- **8x 16-bit analogue in (22.05kHz)**
- **8x 16-bit analogue out (22.05kHz)**

Find an interactive pin out diagram at <http://bela.io/belaDiagram>

Getting Started

bela.io/code/wiki

Materials

what you need to get started...

- **BeagleBone Black** (BBB)
- **Bela Cape**
- **SD card** with Bela image
- 3.5mm headphone jack **adapter cable**
- **Mini-USB cable** (to attach BBB to computer)
- Also useful for hardware hacking: **breadboard**, **jumper wires**, etc.

Step 1

install BBB drivers and Bela software

Install the BeagleBone Black drivers for your OS:

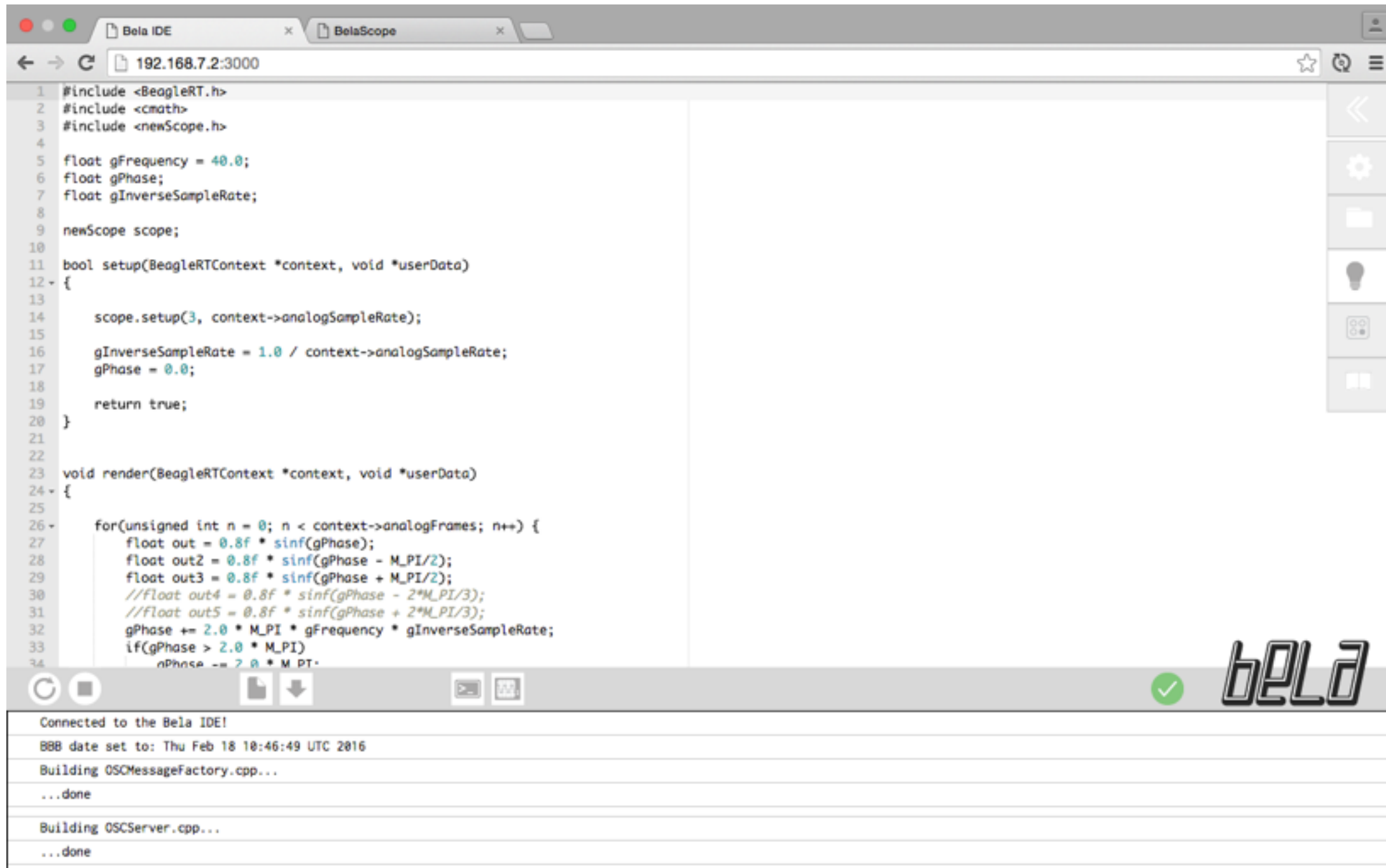
<http://bela.io/code/wiki> --> Getting Started

Bela code (for later):

<http://bela.io/code> --> Downloads --> bela-ableton-workshop.zip

Step 2: Access the IDE:

<http://192.168.7.2:3000>



```
1 #include <BeagleRT.h>
2 #include <cmath>
3 #include <newScope.h>
4
5 float gFrequency = 40.0;
6 float gPhase;
7 float gInverseSampleRate;
8
9 newScope scope;
10
11 bool setup(BeagleRTContext *context, void *userData)
12 {
13     scope.setup(3, context->analogSampleRate);
14
15     gInverseSampleRate = 1.0 / context->analogSampleRate;
16     gPhase = 0.0;
17
18     return true;
19 }
20
21
22
23 void render(BeagleRTContext *context, void *userData)
24 {
25     for(unsigned int n = 0; n < context->analogFrames; n++) {
26         float out = 0.8f * sinf(gPhase);
27         float out2 = 0.8f * sinf(gPhase - M_PI/2);
28         float out3 = 0.8f * sinf(gPhase + M_PI/2);
29         //float out4 = 0.8f * sinf(gPhase - 2*M_PI/3);
30         //float out5 = 0.8f * sinf(gPhase + 2*M_PI/3);
31         gPhase += 2.0 * M_PI * gFrequency * gInverseSampleRate;
32         if(gPhase > 2.0 * M_PI)
33             gPhase -= 2.0 * M_PI;
34     }
```

Connected to the Bela IDE!

BBB date set to: Thu Feb 18 10:46:49 UTC 2016

Building OSCMessageFactory.cpp...

...done

Building OSCServer.cpp...

...done

API introduction

- In `render.cpp`....
- Three main functions:
- `setup()`
runs once at the beginning, before audio starts
gives channel and sample rate info
- `render()`
called repeatedly by Bela system ("callback")
passes input and output buffers for audio and sensors
- `cleanup()`
runs once at end
release any resources you have used
- bela.io/code/embedded *Code docs*

First test program

```
float gPhase; /* Phase of the oscillator (global variable) */  
  
void render(BeagleRTContext *context, void *userData)  
{  
    /* Iterate over the number of audio frames */  
    for(unsigned int n = 0; n < context->audioFrames; n++) {  
        /* Calculate the output sample based on the phase */  
        float out = 0.8f * sinf(gPhase);  
  
        /* Update the phase according to the frequency */  
        gPhase += 2.0 * M_PI * gFrequency * gInverseSampleRate;  
        if(gPhase > 2.0 * M_PI)  
            gPhase -= 2.0 * M_PI;  
  
        /* Store the output in every audio channel */  
        for(unsigned int channel = 0;  
            channel < context->audioChannels; channel++)  
            context->audioOut[n * context->audioChannels  
                + channel] = out;  
    }  
}
```

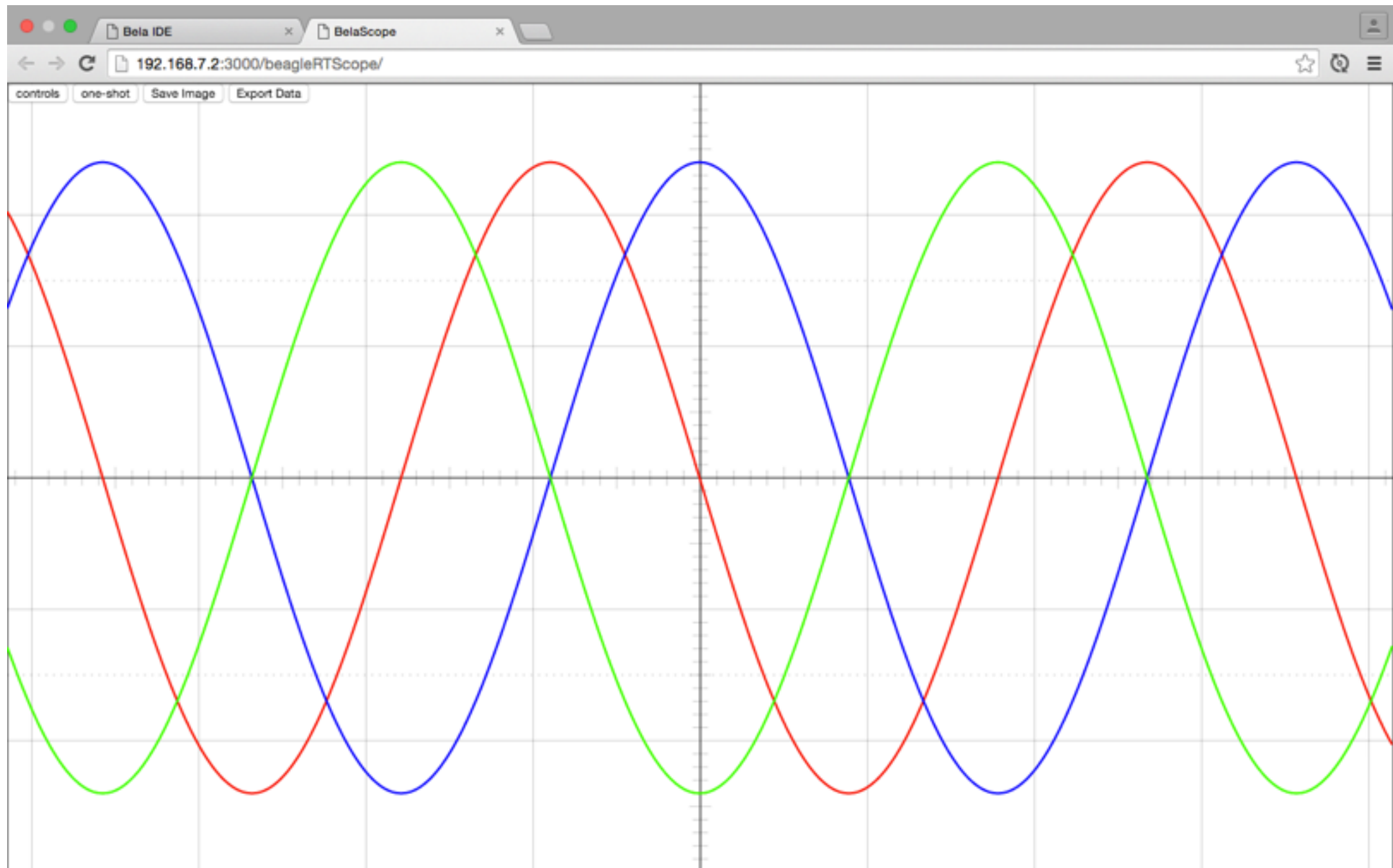
This runs **once per block**

This runs **once per sample** in the block
(**audioFrames** gives the number)

This runs **twice per frame**, once for each channel

One-dimensional array holding interleaved audio data

Access the IDE:
<http://192.168.7.2:3000>



Connect a Potentiometer

a.k.a. a “pot” or knob

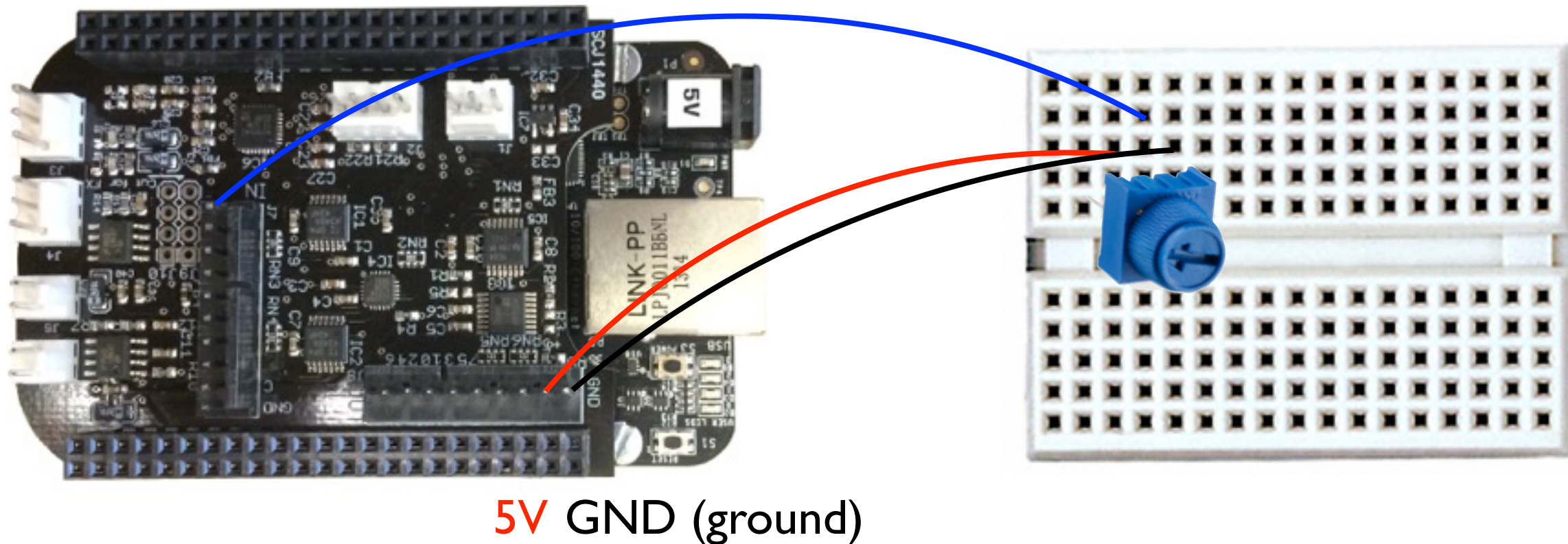
Interactive pinout: <http://bela.io/belaDiagram>

The pot has 3 pins

5V and GND on the outside

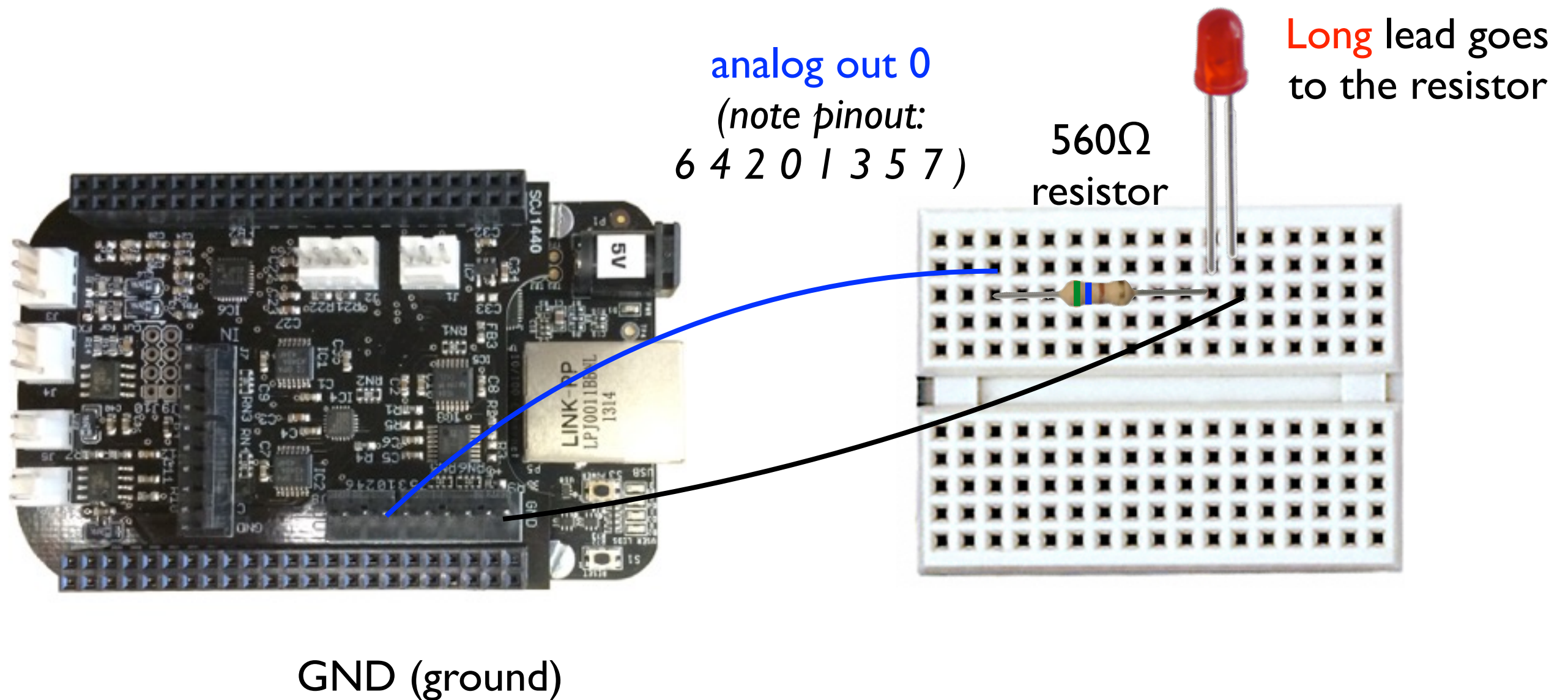
Bela **analog in** in the middle

analog in 0



Connect an LED*

* Light-Emitting Diode



How to build other projects

1. **Web interface:** <http://192.168.7.2:3000>
Edit and compile code on the board
2. **Building scripts:**
 1. **Heavy Pd-to-C compiler** (<https://enzienaudio.com>)
Make audio patches in Pd-vanilla, translate to C and compile on board
 2. **Libpd**
Compile Pd patches without Heavy - access to more objects but not as fast, but good for prototyping
 3. **Faust**
Build online, export to C++, run on Bela

Bela and PureData

Heavy	libpd
Proprietary compiler, cloud-based, MIT non-commercial code	Free
Targets a variety of platforms (C, js, Unity, VST2)	Many ports (ofxpd, webpd ...)
Requires internet connection and local compiling (~1minute)	Instantaneous (save the pd patch and restart)
Generates fast, optimized code, uses little CPU	It is just Pd (...)

libpd on Bela

How to run PureData patches on Bela with libpd :

1. Go to <http://bela.io/code/files> and download the bela-ableton-2016-04-12.zip archive
2. Unzip the archive into a convenient location and open a terminal window
3. Navigate into the scripts/ folder and run
`./run_pd_libpd.sh` `../projects/heavy/pd/demo-track/`
4. Type "yes" and you should hear something

Bela and Faust

- Today: you will have to download the C++ file generated by the <http://faust.grame.fr/onlinecompiler/> (after setting the -i flag), save it on your computer and target it with the build_project.sh script, as in:

```
/path/to/bela/repo/scripts/build_project.sh /path/to/faust/file/CppCode.cpp
```

```
freq = hslider("[1]Frequency[BELA:ANALOG_0]",  
440,460,1500,1):smooth(0.999);  
pressure = hslider("[2]Pressure[style:knob][BELA:ANALOG_4]", 0.96, 0.2,  
2.0, 0.01):smooth(0.999):min(0.99):max(0.2);  
gate = hslider("[0]ON/OFF (ASR Envelope)[BELA:DIGITAL_0]",0,0,1,1);
```

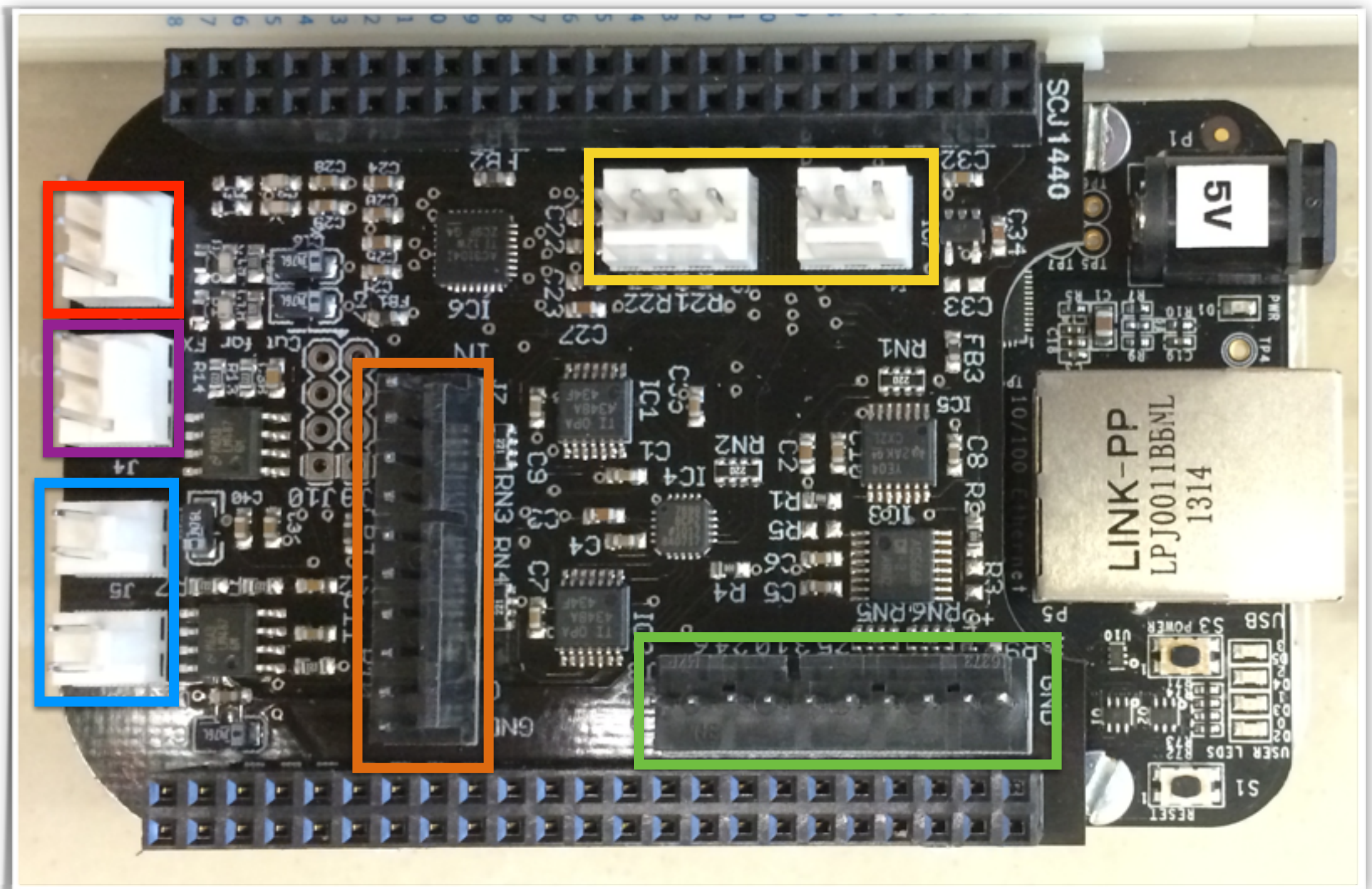
Bela Cape

I2C and GPIO

Audio In

Audio Out
(headphone)

Speakers



8-ch. 16-bit ADC

8-ch. 16-bit DAC

Connect a Potentiometer

a.k.a. a “pot” or knob

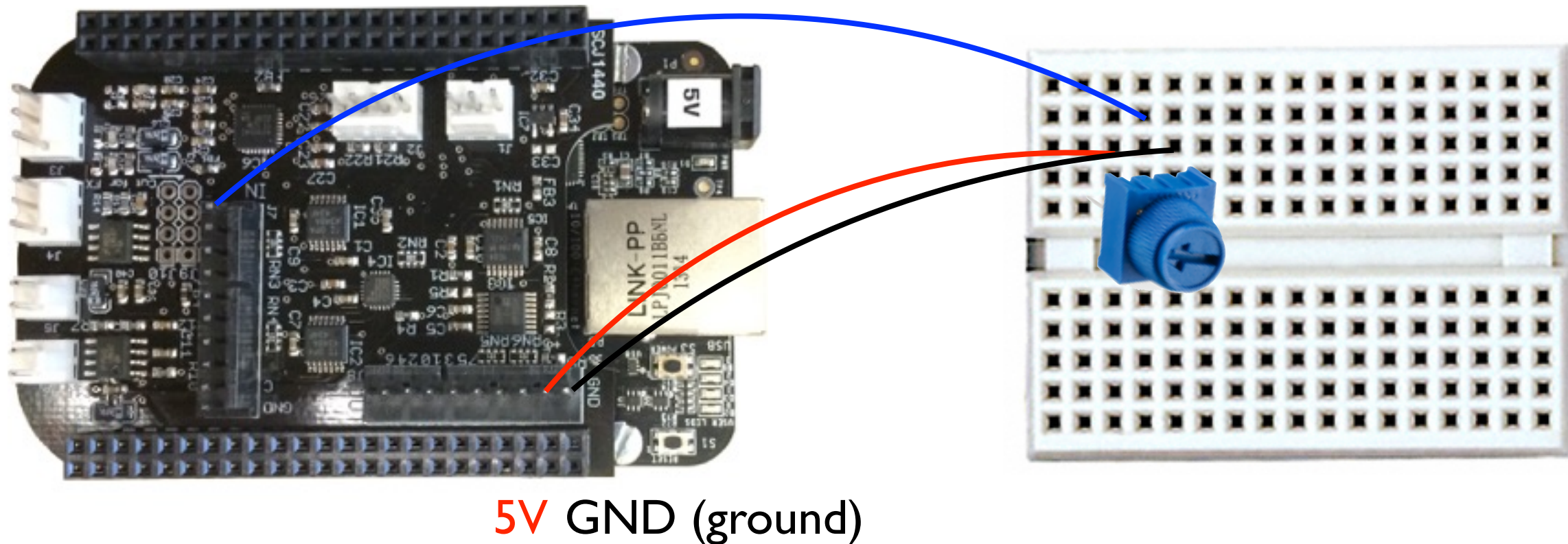
Interactive pinout: <http://bela.io/belaDiagram>

The pot has 3 pins

5V and GND on the outside

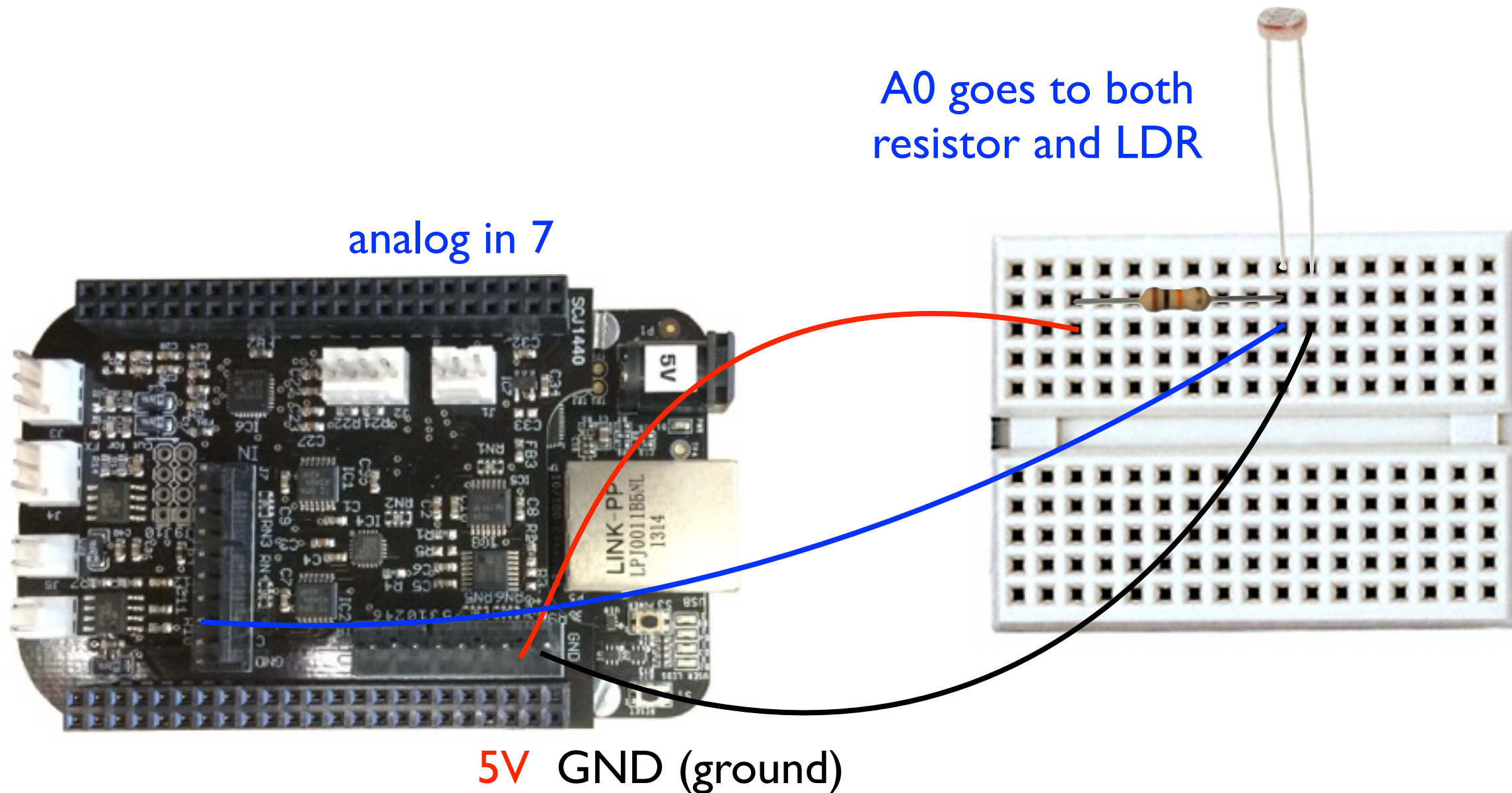
Bela **analog in** in the middle

analog in 0



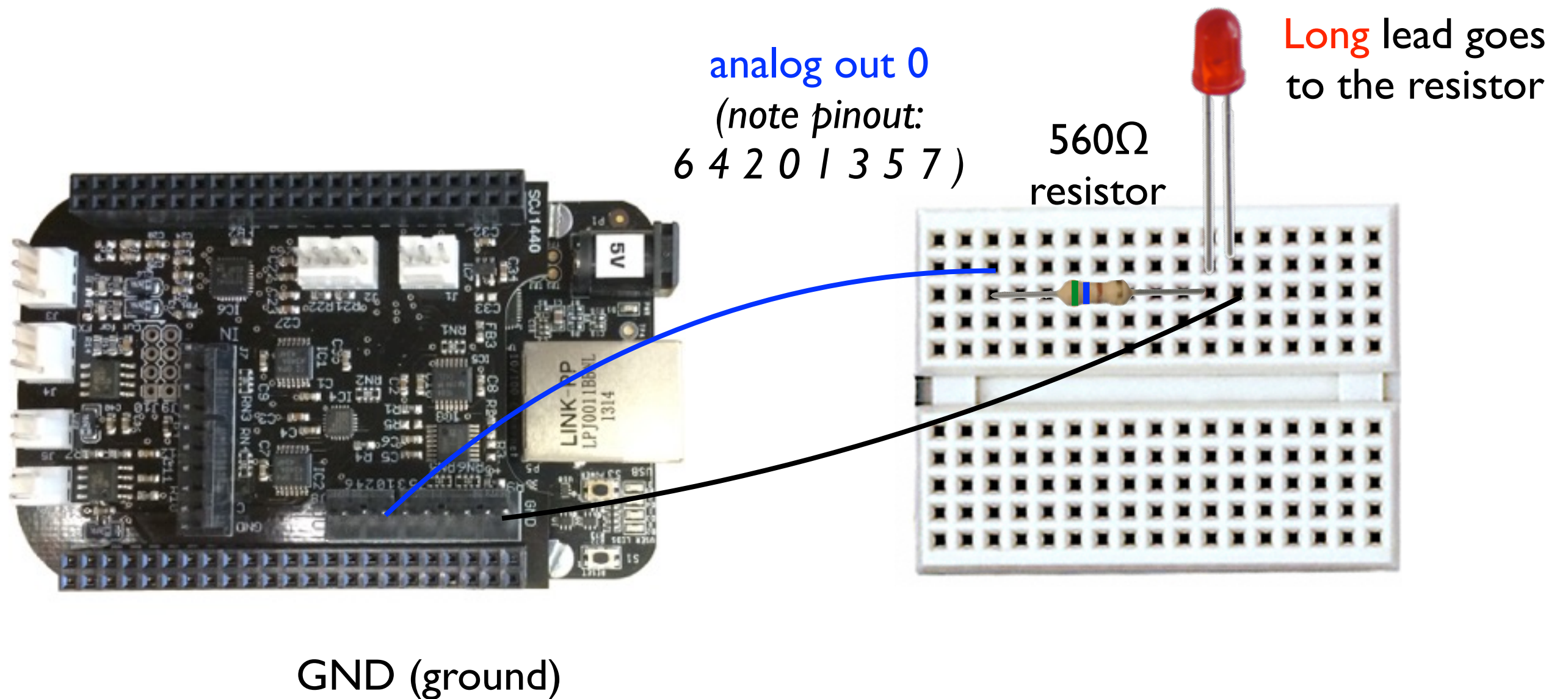
Connect a LDR/FSR*

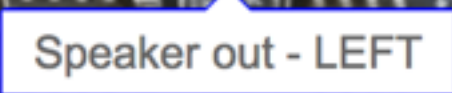
* Light-Dependent Resistor / Force-Sensing Resistor



Connect an LED*

* Light-Emitting Diode





API introduction

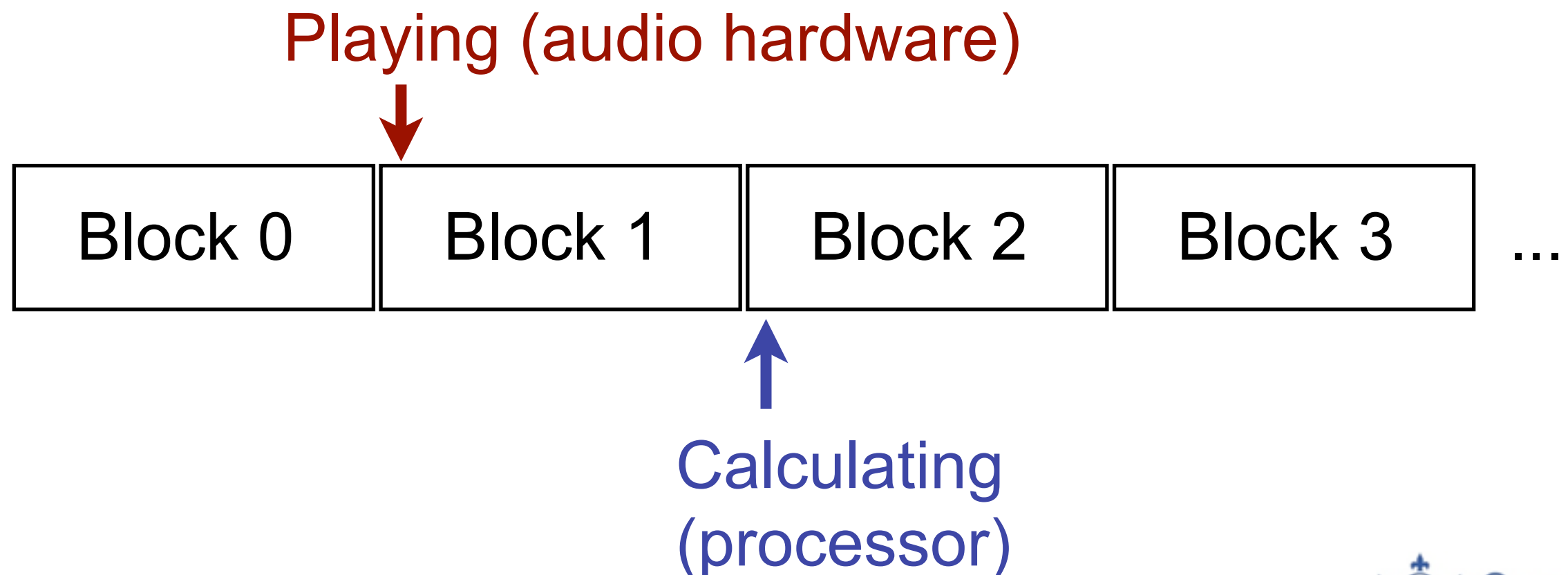
- In `render.cpp`....
- Three main functions:
- `setup()`
runs once at the beginning, before audio starts
gives channel and sample rate info
- `render()`
called repeatedly by Bela system ("callback")
passes input and output buffers for audio and sensors
- `cleanup()`
runs once at end
release any resources you have used

Real-time audio

- Suppose we have code that runs **offline**
 - ▶ (non-real time)
- Our goal is to re-implement it **online** (real time)
 - ▶ Generate audio **as we need it!**
 - ▶ Why couldn't we just generate it all in advance, and then play it when we need it?
- Digital audio is composed of **samples**
 - ▶ 44100 samples per second in our example
 - ▶ That means we need a new sample every $1/44100$ seconds (about every $23\mu\text{s}$)
 - ▶ So option #1 is to run a short bit of code every sample whenever we want to know what to play next
 - ▶ What might be some drawbacks of this approach?
 - Can we guarantee we'll be ready for each new sample?

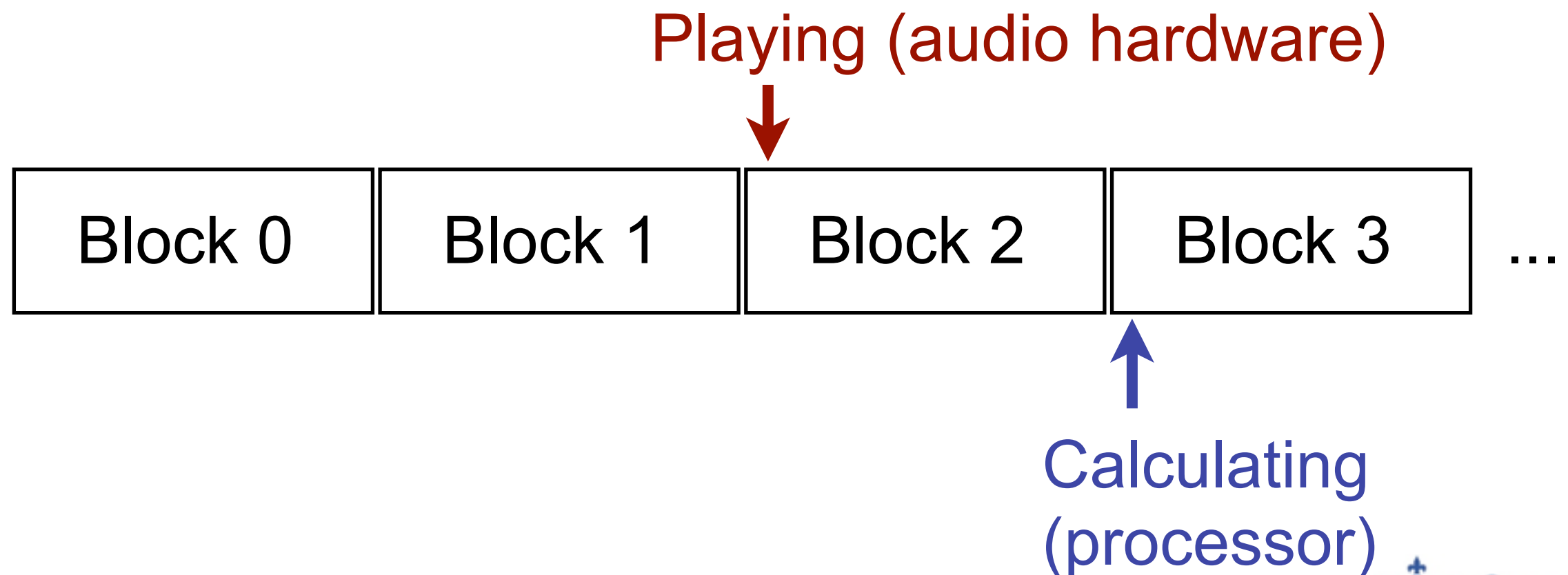
Block-based processing

- Option #2: Process in **blocks** of several samples
 - ▶ Basic idea: generate enough samples to get through the next few milliseconds
 - ▶ Typical **block sizes**: 32 to 1024 samples
 - Usually a power of 2 for reasons having to do with hardware
 - ▶ While the audio hardware is busy playing one block, we can start calculating the next one so it's ready on time:



Block-based processing

- Option #2: Process in **blocks** of several samples
 - ▶ Basic idea: generate enough samples to get through the next few milliseconds
 - ▶ Typical **block sizes**: 32 to 1024 samples
 - Usually a power of 2 for reasons having to do with hardware
 - ▶ While the audio hardware is busy playing one block, we can start calculating the next one so it's ready on time:

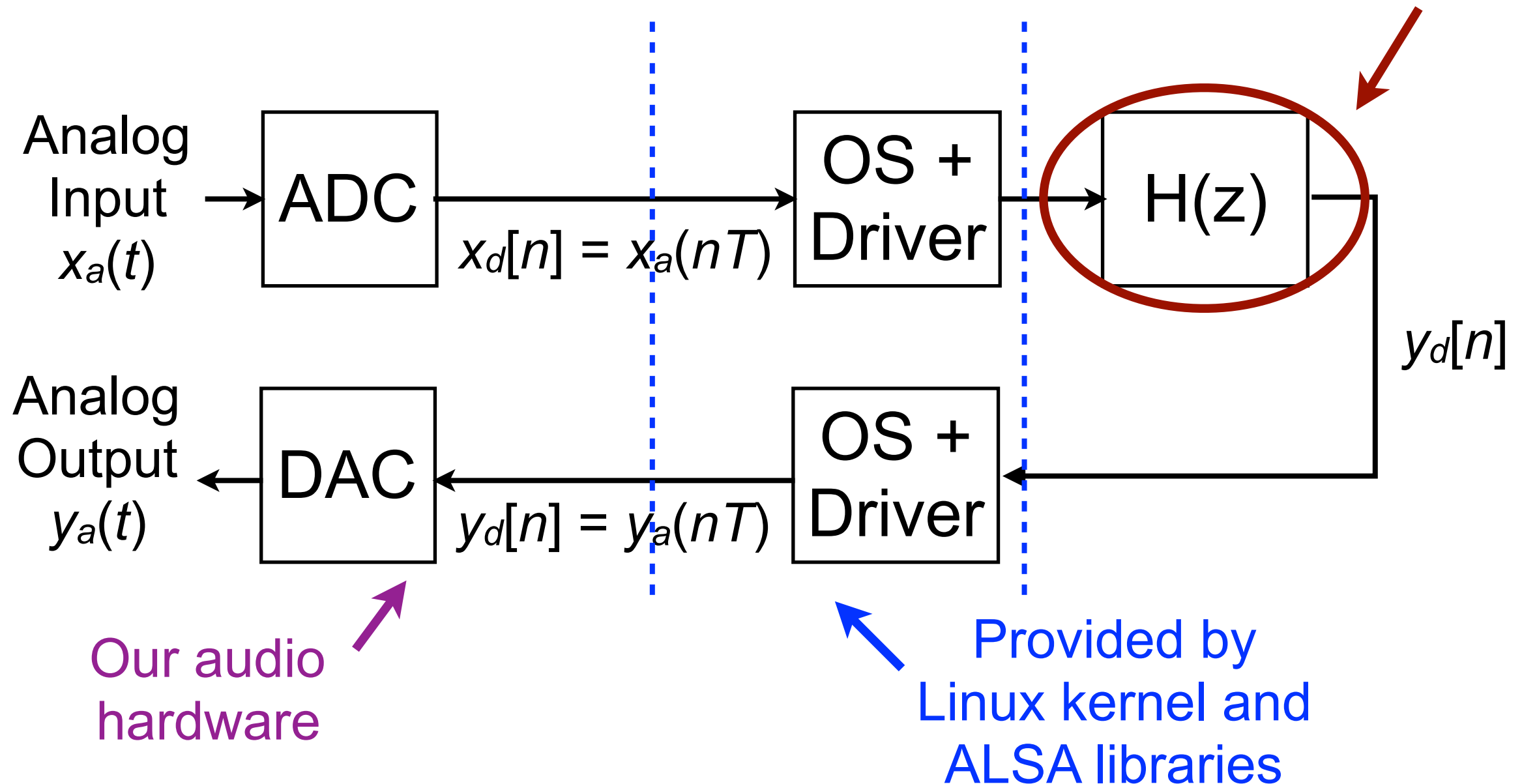


Block-based processing

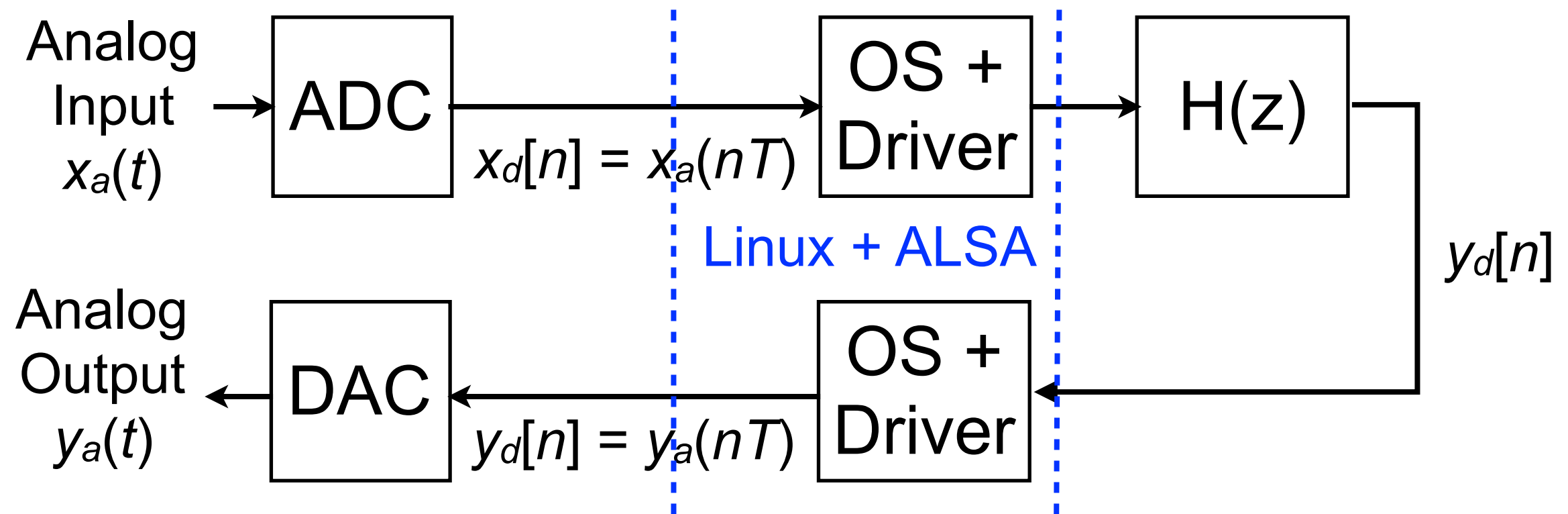
- Advantages of blocks over individual samples
 - ▶ We need to run our function less often
 - ▶ We always **generate one block ahead** of what is actually playing
 - ▶ Suppose one block of samples lasts 5ms, and running our code takes 1ms
 - ▶ Now, we can **tolerate a delay** of up to 4ms if the OS is busy with other tasks
 - ▶ Larger block size = can tolerate more variation in timing
- What is the disadvantage?
 - ▶ **Latency (delay)**

Latency

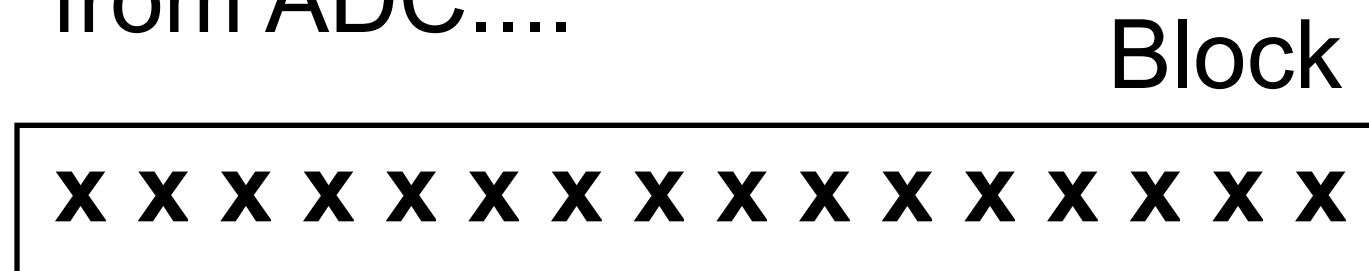
- Primary tradeoff for buffering: **latency**
 - ▶ There will be a **delay** from input to output
- Let's consider a full-duplex system (in and out)
 - ▶ Which are the sources of latency? **We have been writing this**



Latency: the role of buffering

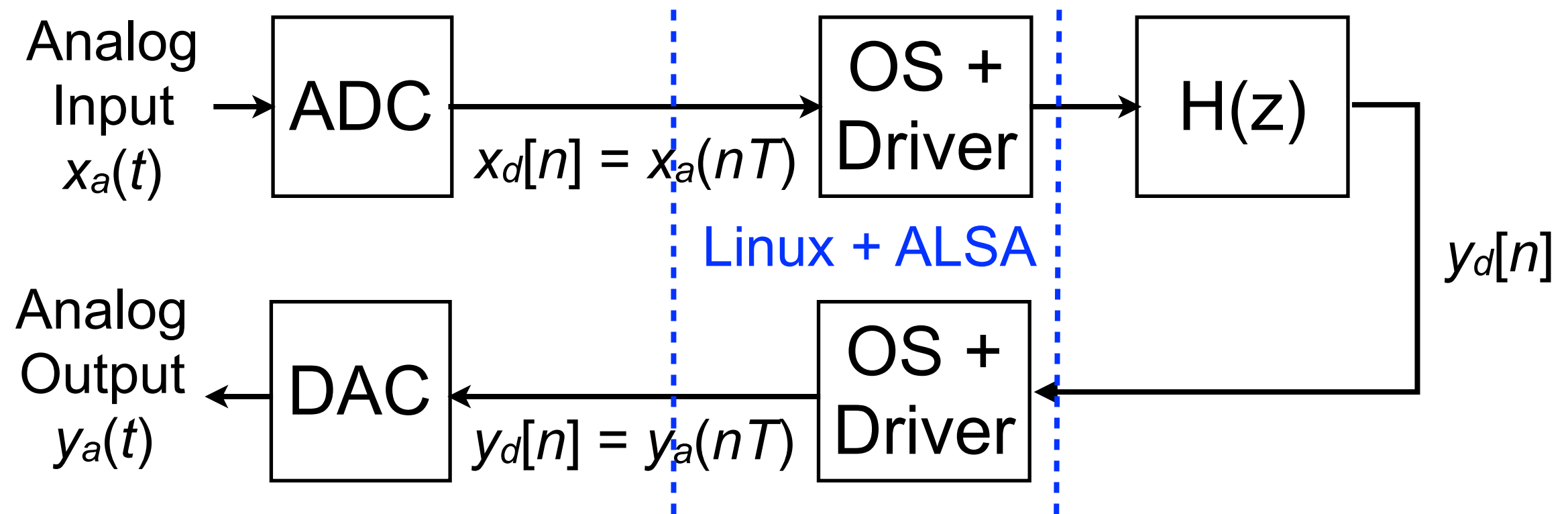


- Block-based processing introduces latency
 - ▶ This is **in addition to** whatever was generated by $H(z)$
- On input side: how is a block of samples created?
 - ▶ For block of size N : we wait until N samples have arrived from ADC....

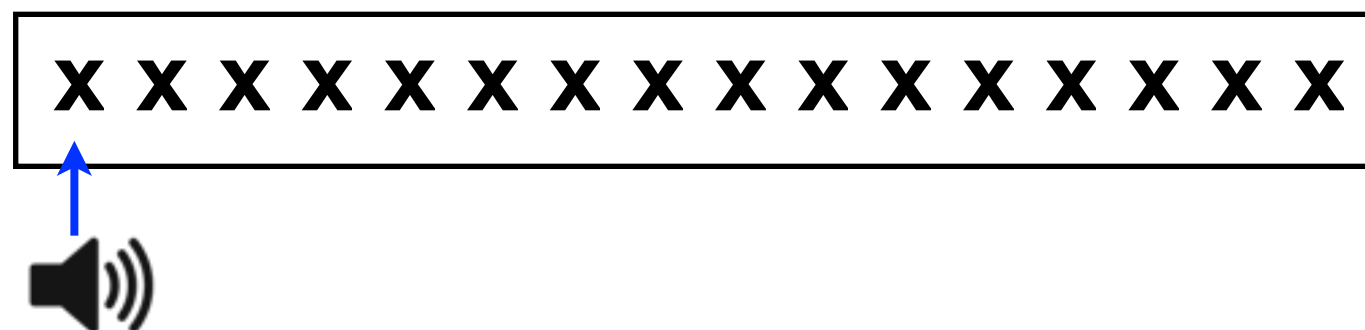


In other words: first sample in the block is already N samples old by the time we get it

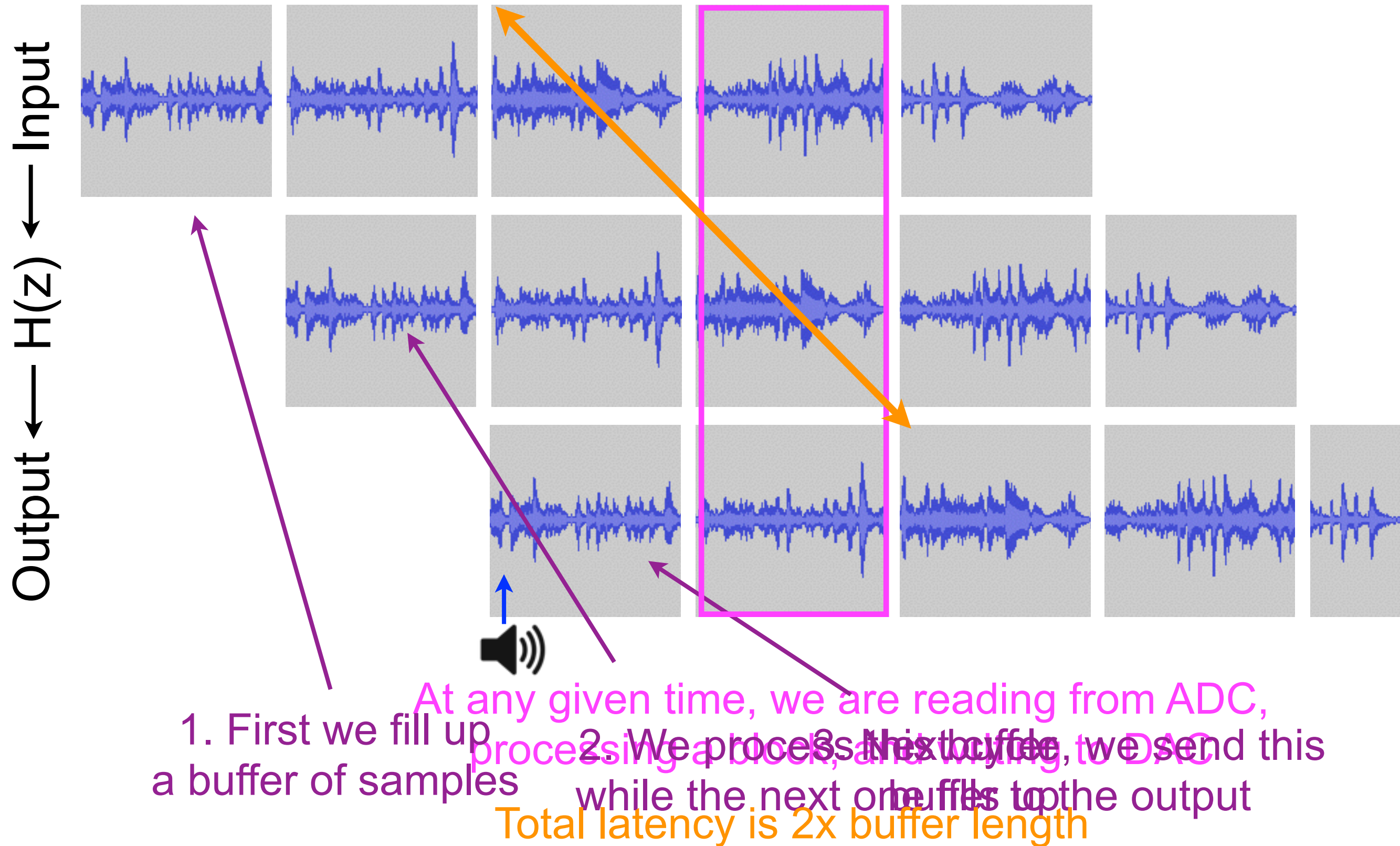
Latency: the role of buffering



- On output side: how is a block played by DAC?
 - ▶ We can only **start** playing once the block arrives!
 - ▶ So how long until the last sample is played?
 - **N samples** after the the block is sent to the hardware



Buffering illustration



API introduction

```
void render(BeagleRTContext *context, void *userData)
```

- Sensor ("matrix" = ADC+DAC) data is gathered **automatically** alongside audio
- Audio runs at **44.1kHz**; sensor data at **22.05kHz**
- **context** holds buffers plus information on number of frames and other info
- Your job as programmer: render one buffer of audio and sensors and finish as soon as possible!
- API documentation: <http://beaglert.cc>

First test program

```
float gPhase; /* Phase of the oscillator (global variable) */

void render(BeagleRTContext *context, void *userData)
{
    /* Iterate over the number of audio frames */
    for(unsigned int n = 0; n < context->audioFrames; n++) {
        /* Calculate the output sample based on the phase */
        float out = 0.8f * sinf(gPhase);

        /* Update the phase according to the frequency */
        gPhase += 2.0 * M_PI * gFrequency * gInverseSampleRate;
        if(gPhase > 2.0 * M_PI)
            gPhase -= 2.0 * M_PI;

        /* Store the output in every audio channel */
        for(unsigned int channel = 0;
            channel < context->audioChannels; channel++)
            context->audioOut[n * context->audioChannels
                + channel] = out;
    }
}
```

This runs **once per block**

This runs **once per sample** in the block
(**audioFrames** gives the number)

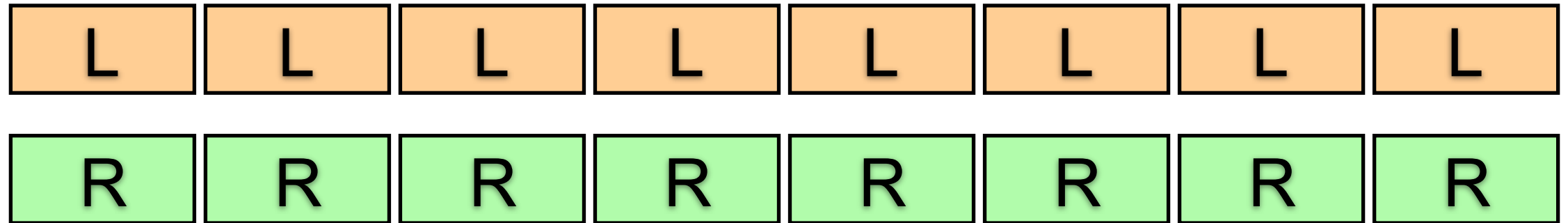
This runs **twice per frame**, once for each channel

One-dimensional array holding interleaved audio data

Interleaving

- Two ways for **multichannel** audio to be stored

- ▶ Way 1: **Separate memory buffers** per channel

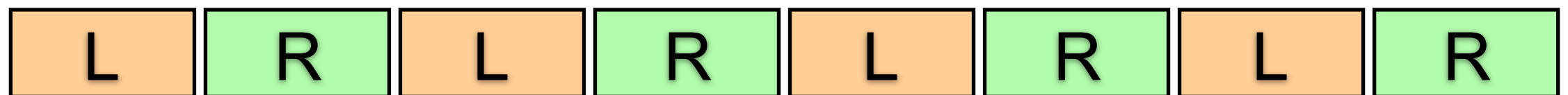


- This is known as **non-interleaved** format
- Typically presented in C as a two-dimensional array:

```
float **sampleBuffers
```

- ▶ Way 2: **One memory buffer** for all channels

- Alternating data between channels



- This is known as **interleaved** format
- Typically presented in C as a one-dimensional array:

```
float *sampleBuffer
```

Interleaving

- We accessed non-interleaved data like this:

- ▶ `float in = sampleBuffers[channel][n];`

- How do we do the same thing with **interleaving**?

- ▶ `float in = sampleBuffers[***what goes here?***];`

- ▶ What else do we need to know?

- Number of channels

1 ch:

L	L	L	L	L	L	L	L
---	---	---	---	---	---	---	---

2 ch:

L	R	L	R	L	R	L	R
---	---	---	---	---	---	---	---

4 ch:

1	2	3	4	1	2	3	4
---	---	---	---	---	---	---	---

- ▶ `float in = sampleBuffers[numChannels*n + channel];`

- ▶ Each sample advances **numChannels** in the buffer

- ▶ The **offset** tells us which channel we're reading

First test program

```
float gPhase; /* Phase of the oscillator (global variable) */

void render(BeagleRTContext *context, void *userData)
{
    /* Iterate over the number of audio frames */
    for(unsigned int n = 0; n < context->audioFrames; n++) {
        /* Calculate the output sample based on the phase */
        float out = 0.8f * sinf(gPhase);

        /* Update the phase according to the frequency */
        gPhase += 2.0 * M_PI * gFrequency * gInverseSampleRate;
        if(gPhase > 2.0 * M_PI)
            gPhase -= 2.0 * M_PI;

        /* Store the output in every audio channel */
        for(unsigned int channel = 0;
            channel < context->audioChannels; channel++)
            context->audioOut[n * context->audioChannels
                + channel] = out;
    }
}
```

This runs **once per block**

This runs **once per sample** in the block
(**audioFrames** gives the number)

This runs **twice per frame**, once for each channel

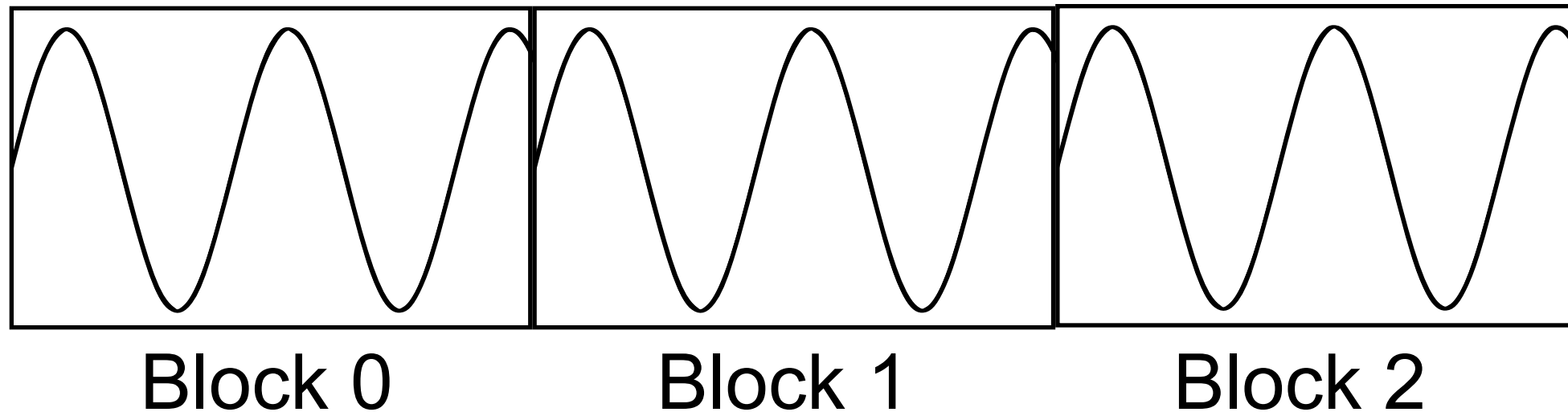
One-dimensional array holding interleaved audio data

Blocks and phase: task

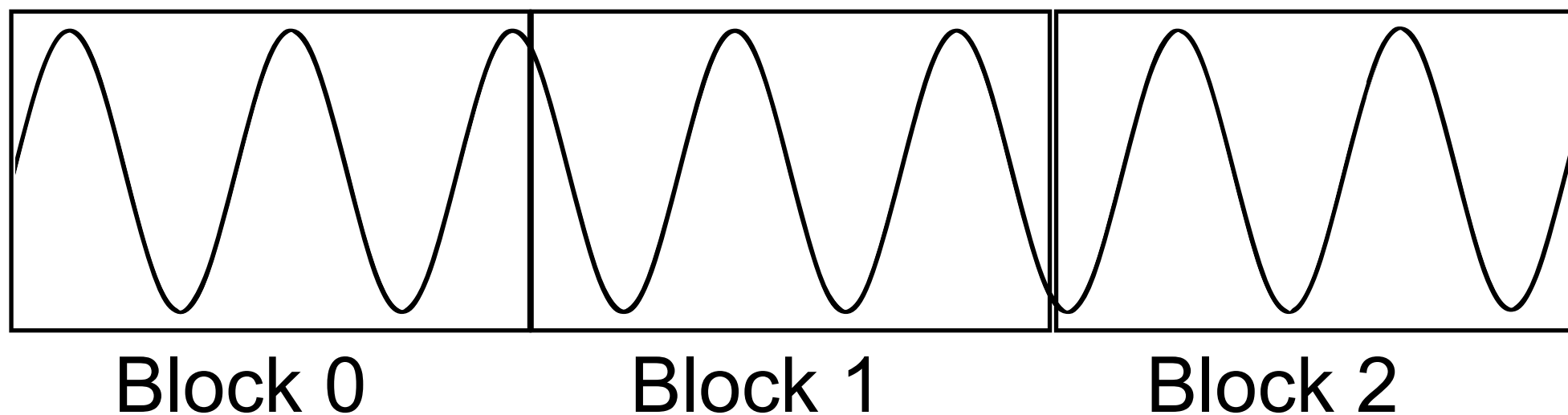
- Need to **preserve state** between calls to `render()`
 - ▶ When you call `render()` a second time, it should **remember where it left off** the first time
 - ▶ But local variables in the function all disappear when the function returns!
 - ▶ Solution: use **global variables** to save the state
 - Okay, cleaner solutions exist: keep a structure that you pass by pointer as an argument to `render()`. Save your state there.
 - Or in C++, use **instance variables** (variables that are declared in the class rather than within a method). But we'll save that for later.

Blocks and phase

- If we don't store phase in a global variable, we get:



- But what we want is this:



First test program

This remembers where we left off

```
float gPhase; /* Phase of the oscillator (global variable) */
```

```
void render(BeagleRTContext *context, void *userData)
{
```

```
    /* Iterate over the number of audio frames */
```

```
    for(unsigned int n = 0; n < context->audioFrames; n++) {
```

```
        /* Calculate the output sample based on the phase */
```

```
        float out = 0.8f * sinf(gPhase);
```

```
        /* Update the phase according to the frequency */
```

```
        gPhase += 2.0 * M_PI * gFrequency * gInverseSampleRate;
```

```
        if(gPhase > 2.0 * M_PI)
```

```
            gPhase -= 2.0 * M_PI;
```

```
        /* Store the output in every audio channel */
```

```
        for(unsigned int channel = 0;
```

```
            channel < context->audioChannels; channel++)
```

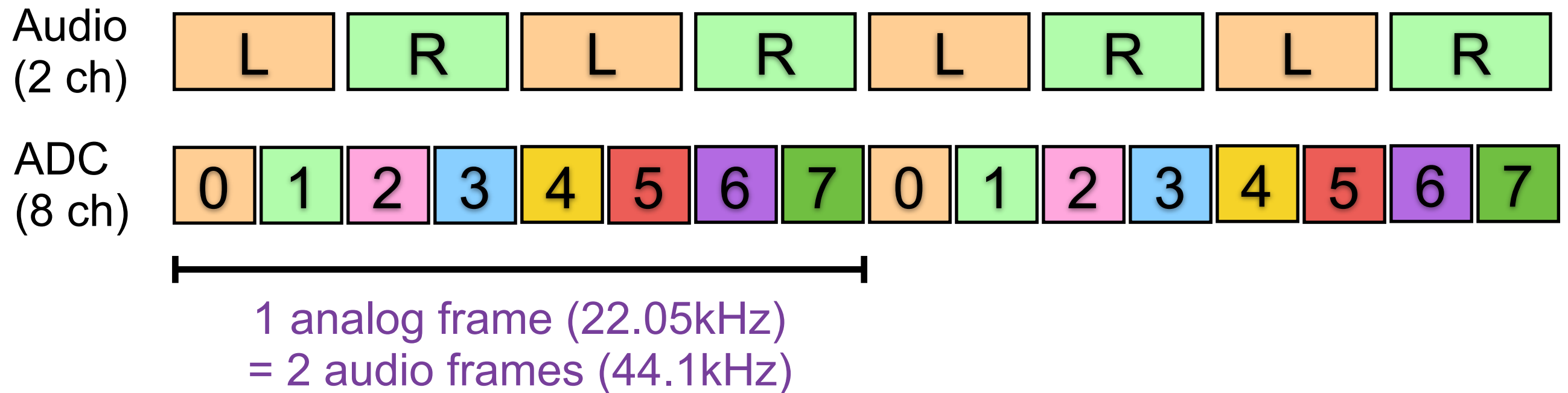
```
            context->audioOut[n * context->audioChannels
                               + channel] = out;
```

```
    }
```

```
}
```

This updates the phase each sample and keeps it in the 0 to 2π range

Analog input data format



- Data type is `float`: just like audio
 - ▶ But range is 0.0 to 1.0
 - This is internally converted from raw values 0 to 65535
 - ▶ Compare this to audio, which is -1.0 to 1.0

Analog input

```
float gPhase;
float gInverseSampleRate;          /* Pre-calculated for convenience */
int gAudioFramesPerAnalogFrame;

extern int gSensorInputFrequency; /* Which analog pin controls frequency */
extern int gSensorInputAmplitude; /* Which analog pin controls amplitude */

void render(BeagleRTContext *context, void *userData)
{
    float frequency = 440.0;
    float amplitude = 0.8;

    for(unsigned int n = 0; n < context->audioFrames; n++) {
        /* There are twice as many audio frames as matrix frames since
           audio sample rate is twice as high */
        if(!(n % gAudioFramesPerAnalogFrame)) {
            /* Every other audio sample: update frequency and amplitude */
            frequency = map(analogReadFrame(context,
                                             n/gAudioFramesPerAnalogFrame,
                                             gSensorInputFrequency),
                           0, 1, 100, 1000);
            amplitude = analogReadFrame(context,
                                         n/gAudioFramesPerAnalogFrame,
                                         gSensorInputAmplitude);
        }

        float out = amplitude * sinf(gPhase);

        for(unsigned int channel = 0; channel < context->audioChannels; channel++)
            context->audioOut[n * context->audioChannels + channel] = out;

        gPhase += 2.0 * M_PI * frequency * gInverseSampleRate;
        if(gPhase > 2.0 * M_PI)
            gPhase -= 2.0 * M_PI;
    }
}
```

This runs **every other sample**

Read the **analog input** at the specified **frame**

Map the 0-1 input range to a frequency range

Digital I/O

```
void render(BeagleRTContext *context, void *userData)
{
    static int count = 0; // counts elapsed samples
    float interval = 0.5; // how often to toggle the LED (in seconds)
    static int status = GPIO_LOW;

    for(unsigned int n = 0; n < context->digitalFrames; n++) {
        /* Check if enough samples have elapsed that it's time to
           blink to the LED */
        if(count == context->digitalSampleRate * interval) {
            count = 0; // reset the counter
            if(status == GPIO_LOW) {
                /* Toggle the LED */
                digitalWriteFrame(context, n, P8_07, status);
                status = GPIO_HIGH;
            }
            else {
                /* Toggle the LED */
                digitalWriteFrame(context, n, P8_07, status);
                status = GPIO_LOW;
            }
        }

        /* Increment the count once per frame */
        count++;
    }
}
```

This runs **once**
per digital frame

Write the **digital**
output at the
specified **frame**

To manage timing, **count**
samples rather than
using delays

bela

Stay tuned! Join the announcement list at

<http://bela.io>