

A Real-time, Spiking SoC for Edge RF Inference

Jon Cowart, Zephan Enciso, Mark Horeni, and Clemens Schaefer

2023-05-14

Contents

1	Overview	2
1.1	System Architecture	2
1.2	Timing and Dataflow	2
2	Case for Acceleration	3
3	Spikifier	3
3.1	Operation	3
3.2	Noise Modeling	4
4	Digital Accelerator	4
4.1	Processing Element Architecture	4
4.2	Physical Implementation	5
5	Division of Labor	5

1 Overview

We present a **real-time, spiking** SoC for **edge radio frequency (RF) inference** in 22 nm CMOS.

1.1 System Architecture

Figure 1 depicts the primary components of the system, which are:

1. **4x4 phased array**, which is mixed down off-chip.
2. **“Spikifier”** analog circuit, which integrates the mixed RF signal during the negative phase of a 10 MHz clock and produces an output spike if the integrated signal surpassed a threshold.
3. **Data buffers** (FIFO) for synchronization.
4. Two layers of **processing elements (PEs)**: 16 layer 1 PEs, which feed 8 layer 2 PEs.
5. **RISC-V processor**, responsible for controlling PEs via a single multicast reset signal and measuring the firing rate of the layer 2 PEs to produce a classification.

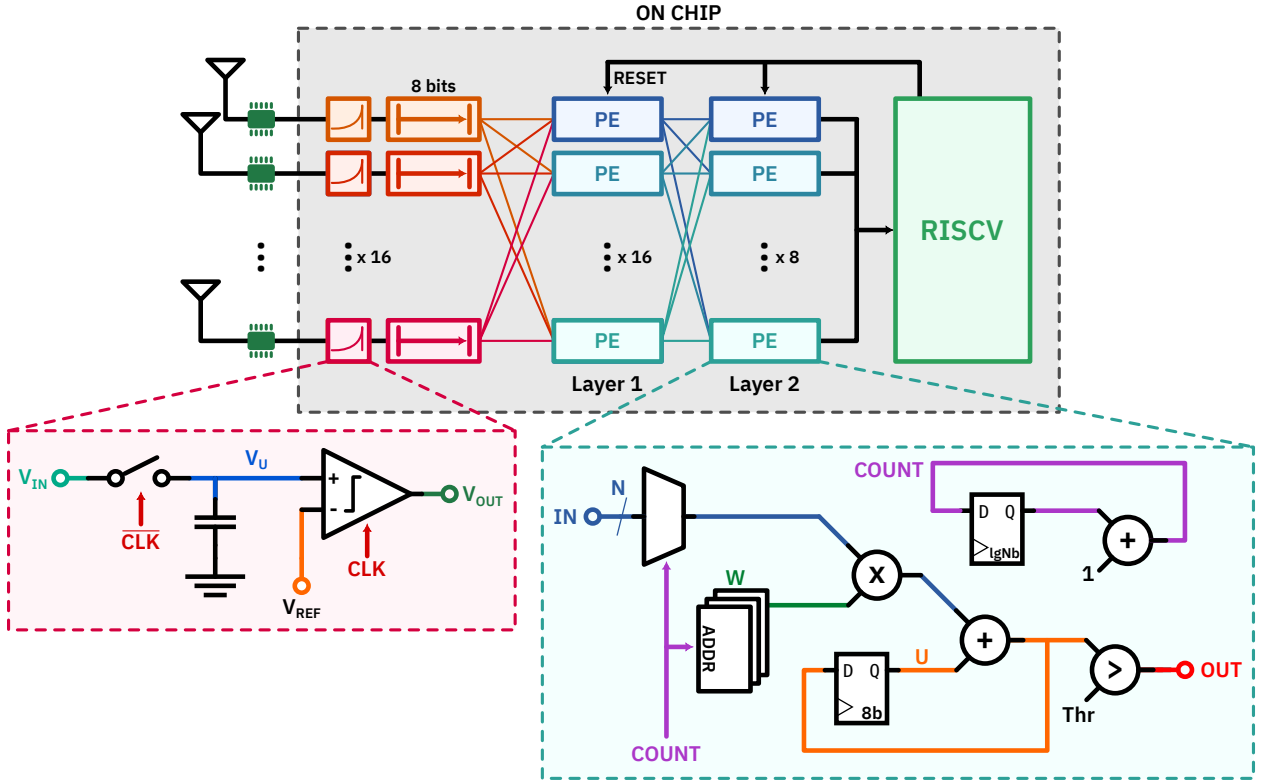


Figure 1: System architecture. Everything within the grey rectangle is on-chip. The red call-out depicts the spikifier circuit, while the teal call-out depicts the PE architecture.

1.2 Timing and Dataflow

The goal of this project is to perform **real-time** inference, which necessitates that the data consumption rate equals the data production rate. The system is fully pipelined (see Figure 2), and the longest single processing step is the layer 1 PEs at 16 clock cycles (the operation of the PE is described in more depth in Section 4.1). Therefore, the core clock must run at 16 times the rate of the sample clock.

The spikifiers integrate the mixed RF signal during the negative phase of a 10 MHz clock and produce an output at the positive edge. FIFOs between the spikifiers and the layer 1 PEs provide synchronization between the 10 MHz sample clock and the 160 MHz core clock. The layer 2 PEs only require 8 clock cycles

to complete their computation, after which the RISC-V processor requires 12 cycles to update its rate counts. The layer 2 PEs may be clock gated during the 8 cycles in which they have nothing to process.

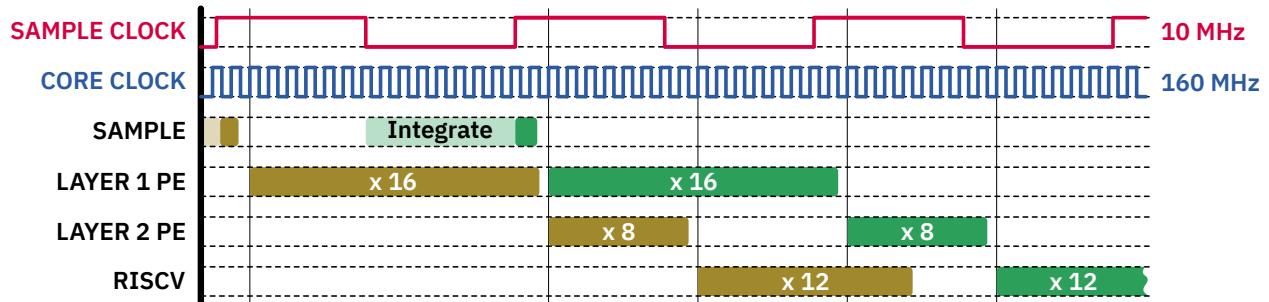


Figure 2: Pipelined system timing. Note the overlap between the spikifier (“sample”), layer 1 PEs, layer 2 PEs, and the RISC-V host processor.

2 Case for Acceleration

For a baseline, we simulate our system on our simple RISC-V core (no out-of-order execution, simple sequential pipelining) using `riscv64-unknown-elf-gcc` from **RISCV-tools** and the `riscv64-unknown-elf` kernel on **Spike**, an ISA simulator. Not including rate counting to produce an output classification, the core computation kernel requires over 180 thousand operations per timestep on the RISC-V core, while our accelerator needs only 16 clock cycles. With real-time processing and an ideal CPI of 1, the un-accelerated version running at the same core clock would only be able to handle an input signal sampled at less than 1 kHz.

3 Spikifier

3.1 Operation

Each spikifier consists of an integration capacitor and a comparator (see Figure 1). The spikifier integrates the input signal onto the capacitor during the negative phase of the clock. If the voltage on the capacitor passes the threshold level, the spikifier will issue an output spike on the positive edge of the clock. The operation is shown below in Figure 3.

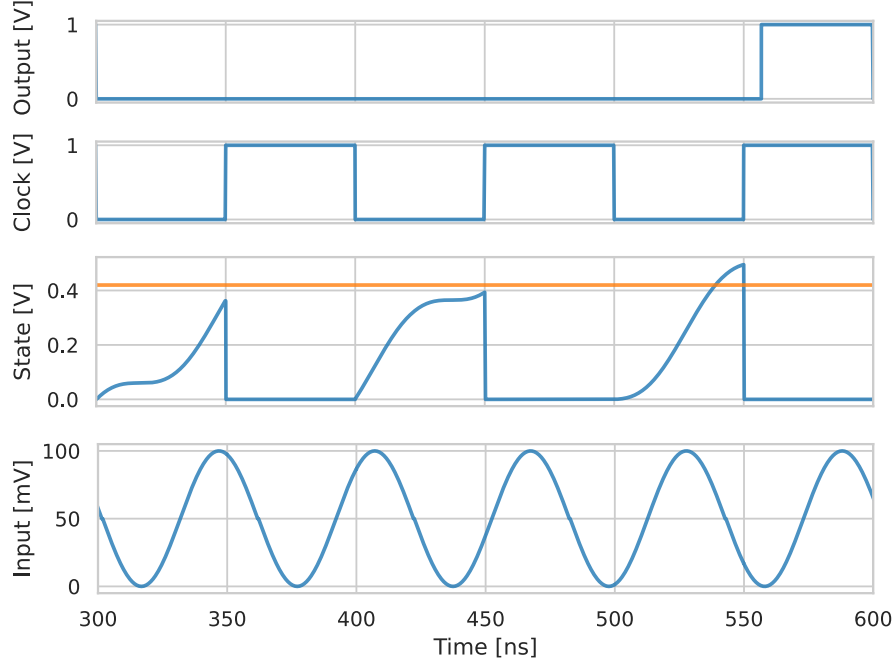


Figure 3: Spikifier operation. The orange line represents the threshold voltage after which the spikifier will output a spike. Since the input in this example is a slow sine wave sampled at 3.4 GHz, the maximum integrated voltage depends on the phase of the input wave when the integration period begins. On the third integration phase (500 to 550 ns), the capacitor voltage (“State”) rises above the threshold, causing an output spike after the rising clock cycle at 550 ns.

3.2 Noise Modeling

There are three principle noise sources in the behavioral model:

1. **Capacitor noise.** The value integrated onto the capacitor experiences additional white noise.
2. **Comparator delay.** The delay between the rising clock edge and the comparator output is modeled as a Gaussian distribution with a mean of 2 ns.
3. **Comparator input-referred noise.** The input to comparator experiences additional white noise.

4 Digital Accelerator

4.1 Processing Element Architecture

Each PE is very simple to maximize density and energy efficiency (see Fig. 1). Each PE has an associated parameter N , which represents the fan-in to that PE.

There are four memory components to the PE: the **membrane potential register** (2 bits), the **weight register file** (N 8-bit registers), the **threshold register**, and the **count register** ($\log N$ bits). The weight register file and threshold register are both filled once on startup through a scan chain. The count and membrane potential register are susceptible to the reset signal issued by the RISC-V processor, but the weight register file and threshold register are not.

With each clock cycle, the count register increments, selecting a new weight from the register file and a new input. If there was a spike on the selected input, then the weight adds to the membrane potential. When the membrane potential exceeds the value stored in the threshold register, the PE produces an output spike.

There are two circumstances which might reset the membrane potential: **rollover** and **reset**. Since the membrane potential is represented with a unsigned integer, exceeding the maximum representable value will cause rollover. As previously mentioned, the membrane potential register is susceptible to the RISC-V reset signal, so it may also be reset externally.

4.2 Physical Implementation

5 Division of Labor

Group Member	Task
Jon	RTL and Validation
Zephan	System Architecture and Analog Behavioral Model
Mark	System Architecture and RISC-V computation kernel
Clemens	Synthesis and Implementation