# SAFARI

# Accelerating Approximate Pattern Matching with Processing-In-Memory (PIM) and Single-Instruction Multiple-Data (SIMD) Programming

Damla Senol Cali[1]

Zülal Bingöl[2], Jeremie S. Kim[1,3],
Rachata Ausavarungnirun[1], Saugata Ghose[1],
Can Alkan[2] and Onur Mutlu[3,1]

[1] **Carnegie Mellon** [2] **Bilkent University** [3] **ETH** *zürich*

# Poster SEQ-15

**Problem:**

- Bitap algorithm can perform string matching with **fast and simple bitwise operations.**
- Due to the **von Neumann architecture**, **memory bus** between the CPU and the memory is the bottleneck of the **high-throughput parallel *bitap* computations**.



**CPU** ← → **Memory**

Low capacity
memory bus

**Goals:**

1) Overcoming **memory bus bottleneck** of approximate string matching by performing **processing-in-memory** to exploit the **high internal bandwidth** available inside new and emerging memory technologies, and

2) Using **SIMD programming** with Xeon Phi to take advantage of the **high amount of bit parallelism** available in the *bitap* algorithm.