# C2HLSC

## Leveraging LLMs to refactor C code into HLS-compatible C

Luca Collini, **Andrew Hennessee**, Ramesh Karri, Siddharth Garg

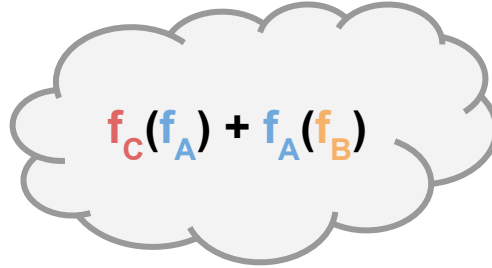NYU | TANDON SCHOOL OF ENGINEERING

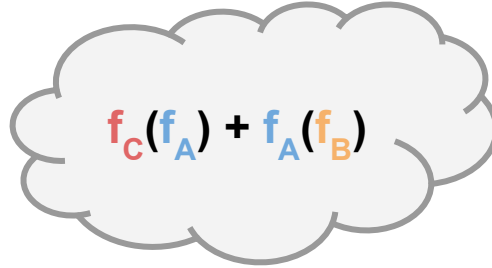# What do system architects do?

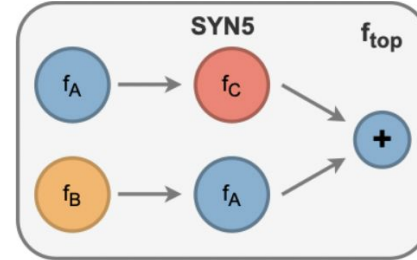# What do system architects do?

**1** Profile application

$$f_C(f_A) + f_A(f_B)$$

# What do system architects do?

**1** Profile application

$$f_C(f_A) + f_A(f_B)$$

**2** Construct DFG



4

# What do system architects do?



**1** Profile application

$$f_C(f_A) + f_A(f_B)$$

**2** Construct DFG

SYN5    $f_{top}$

$f_A \rightarrow f_C$

$f_B \rightarrow f_A$

$+$

**3** Model latency    $f_{top} + \max(f_A + f_C, f_B + f_A)$

# What do system architects do?

1. Profile application

$$f_C(f_A) + f_A(f_B)$$

2. Construct DFG

SYN5

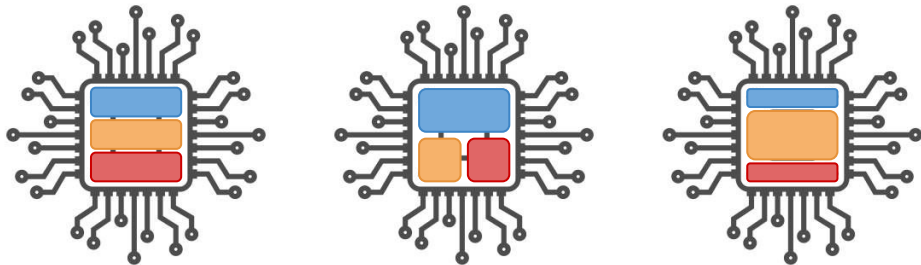$f_A \rightarrow f_C$

$f_B \rightarrow f_A$

$f_{top}$ +

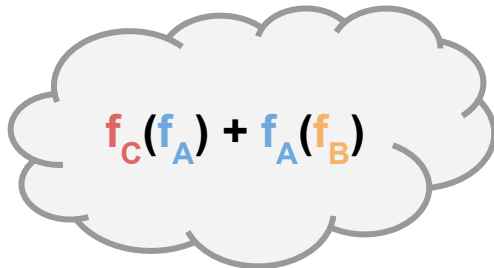3. Model latency

$$f_{\text{top}} + \max(f_A + f_C, f_B + f_A)$$

4. DSE to map kernels to hardware

# Can LLMs replace system architects?
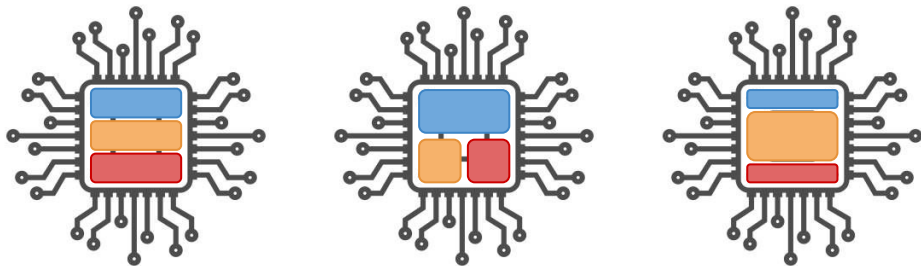
**1** Profile application

$$f_C(f_A) + f_A(f_B)$$

**2** Construct DFG

SYN5

$f_A \rightarrow f_C$

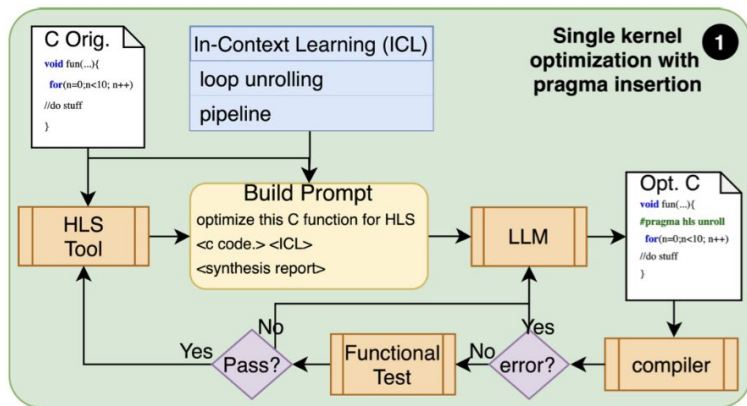$f_B \rightarrow f_A$

$+$ $f_{top}$

**3** Model latency $\quad f_{\text{top}} + \max(f_A + f_C, f_B + f_A)$

**4** DSE to map kernels to hardware

# Overview of proposed flow



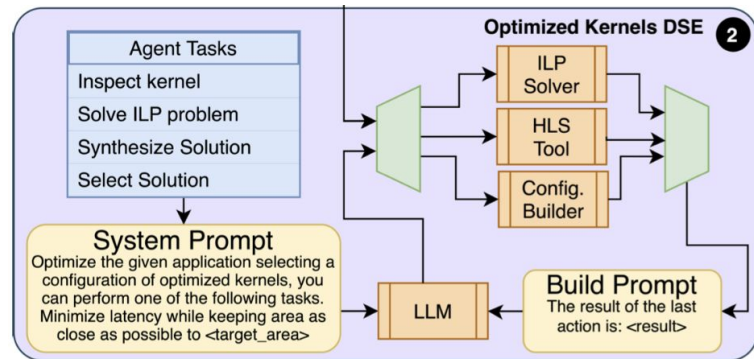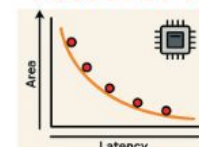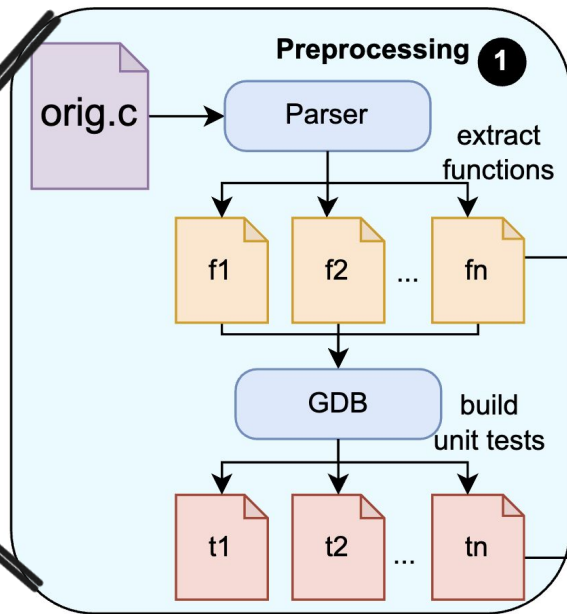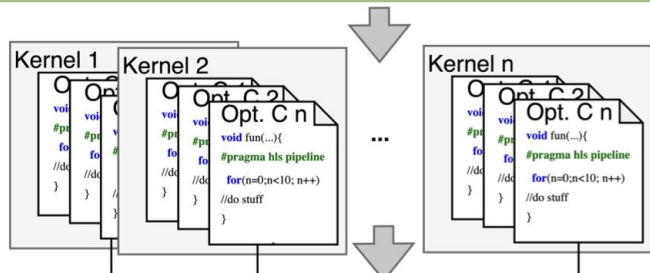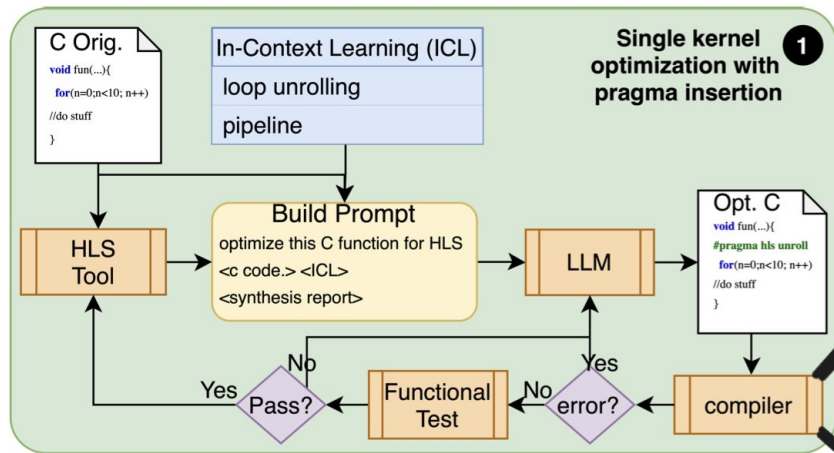**Single Kernel Optimization**

**Full System Composition**

# Single kernel optimization via pragma insertion



(C2HLSC, TODAES)

# Design points for DSE



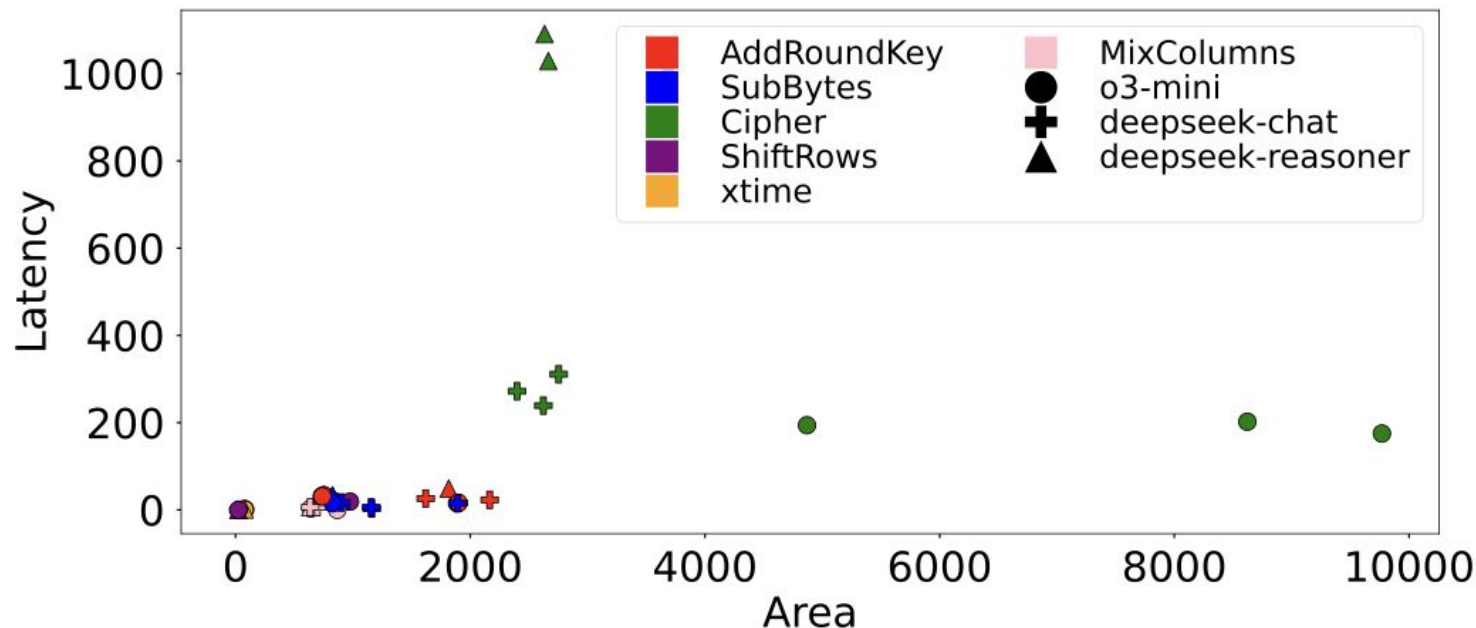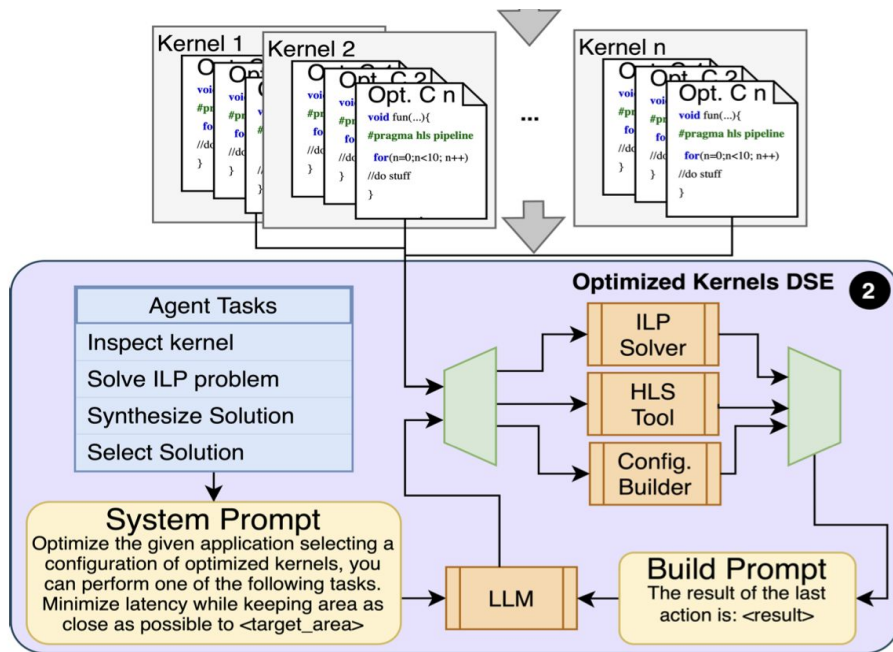Fig. 4: Solutions for AES sub-kernels for each model.

# DSE of optimized kernels

# Thank you!