

# Timing-Driven Placement for FPGA Architectures with Dedicated Routing Paths

---



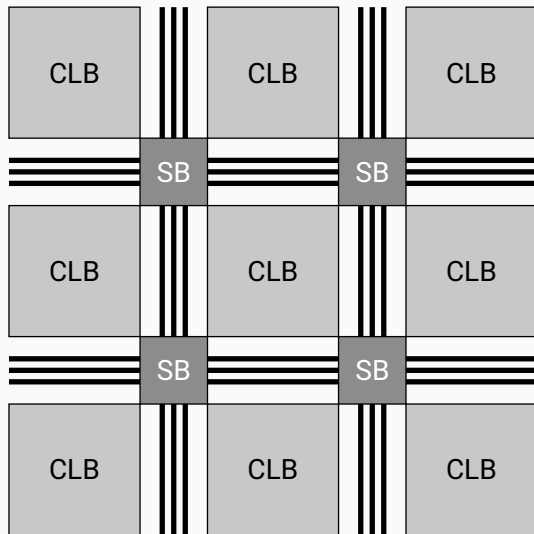
S. Nikolić, G. Zgheib\*, and P. lenne

FPL'20, Göteborg, 01.09.2020

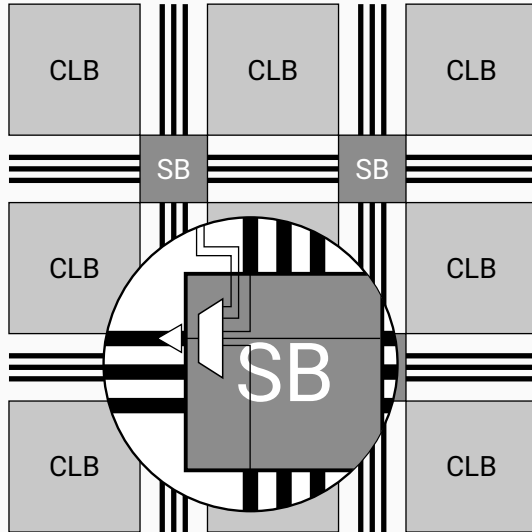
École Polytechnique Fédérale de Lausanne

\*Intel Corporation

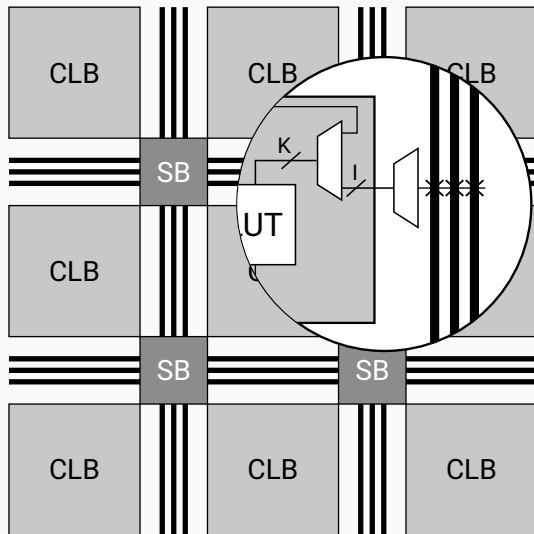
# Field-Programmable Gate Array



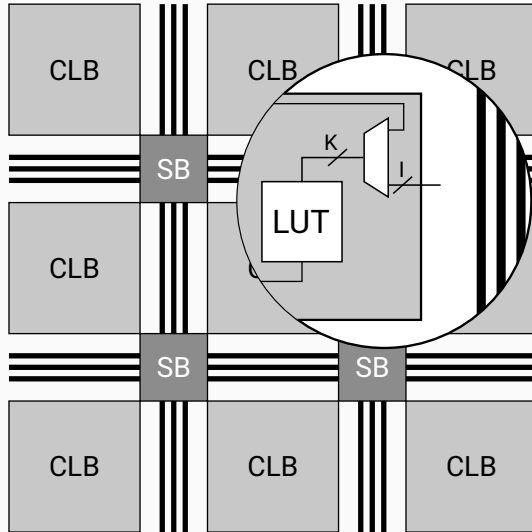
# Price of Programmability: Switch Block MUX



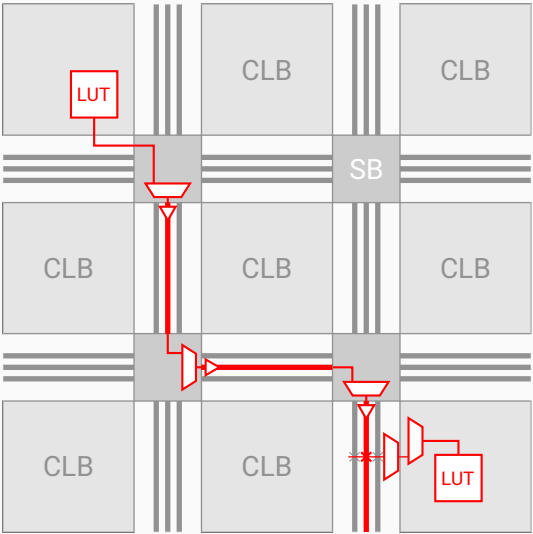
# Price of Programmability: Connection Block MUX




# Price of Programmability: Crossbar MUX



# Many MUXes $\implies$ Large Delay



# Direct Connections: Switch Block-to-Switch Block

Virtex-II 1.5V Field-Programmable Gate Arrays 

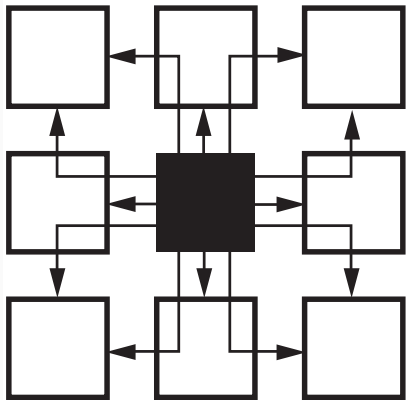
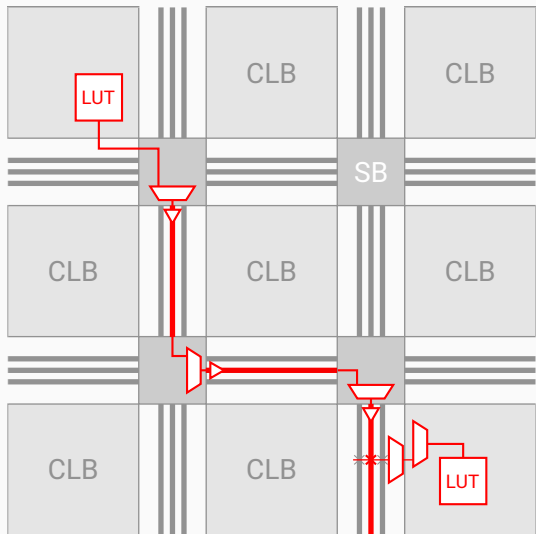


Figure 48: Hierarchical Routing Resources



Module 2 of 4 [www.xilinx.com](http://www.xilinx.com) DS031-2 (v1.9) November 29, 2001  
78 1-800-255-7778 Advance Product Specification

# Direct Connections: Switch Block-to-Switch Block

Virtex-II 1.5V Field-Programmable Gate Arrays 

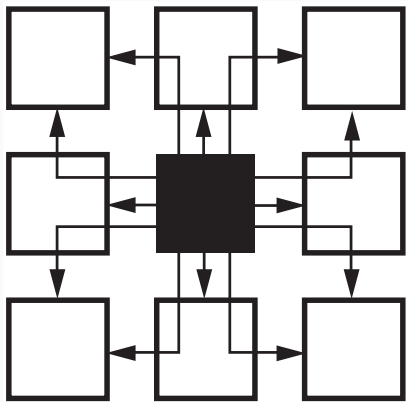
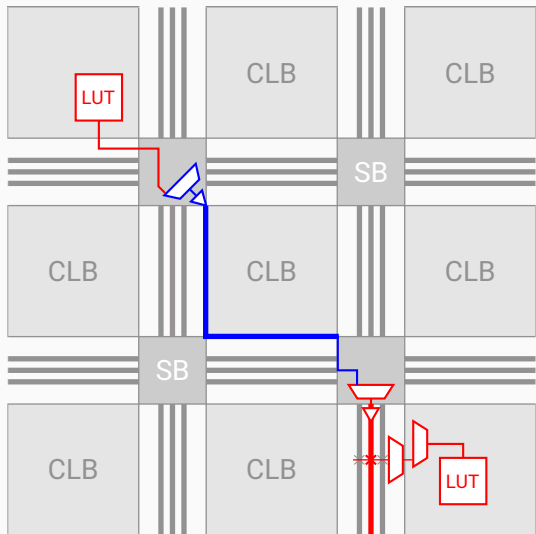


Figure 48: Hierarchical Routing Resources

Module 2 of 4 [www.xilinx.com](http://www.xilinx.com) DS031-2 (v1.9) November 29, 2001  
78 1-800-255-7778 Advance Product Specification





# Direct Connections: Cluster-to-Cluster

UG-510LAB | 2020\_04\_24

 Send Feedback

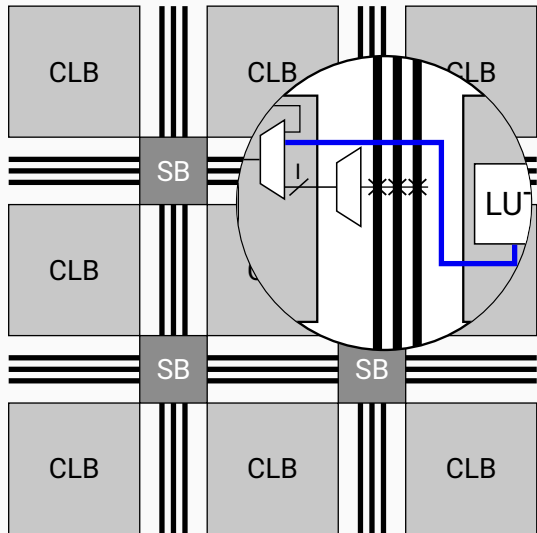
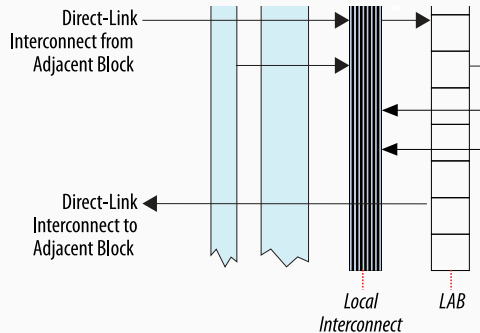


Figure 1. Intel Stratix 10 LAB Structure and Interconnects Overview

# Direct Connections: LUT-to-LUT

## Xilinx Adaptive Compute Acceleration Platform: Versal™ Architecture

Brian Gaide, Dinesh Gaitonde, Chirag Ravishankar, Trevor Bauer  
bgaide@xilinx.com, dineshg@xilinx.com, chiragr@xilinx.com, trevor@xilinx.com  
Xilinx Inc.

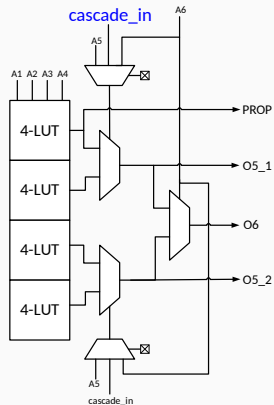
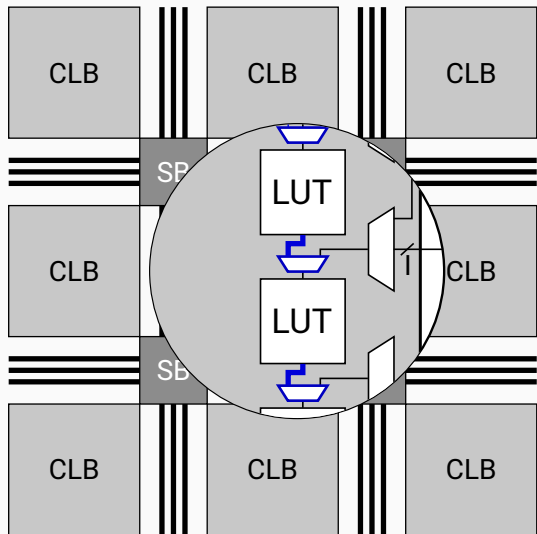


Figure 5: 6LUT Comparison between UltraScale and Versal



# Direct Connections: Two Questions

1. Where to put them?  
(metal and area cost, increased capacitive loading, etc.)

# Direct Connections: Two Questions

1. Where to put them?  
(metal and area cost, increased capacitive loading, etc.)
  
2. How to use them effectively?

# Direct Connections: Two Questions

1. Where to put them?  
(metal and area cost, increased capacitive loading, etc.)

## Our work at FPGA'20

2. How to use them effectively?

# Direct Connections: Two Questions

1. Where to put them?  
(metal and area cost, increased capacitive loading, etc.)

Our work at FPGA'20

2. How to use them effectively?

This work

# Outline

Introduction

Target Architectures

General Approach

Placement Algorithm

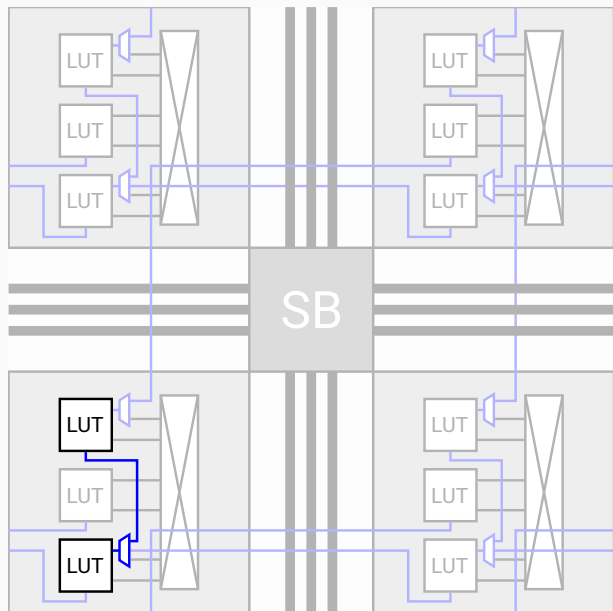
Results

# Target Architectures

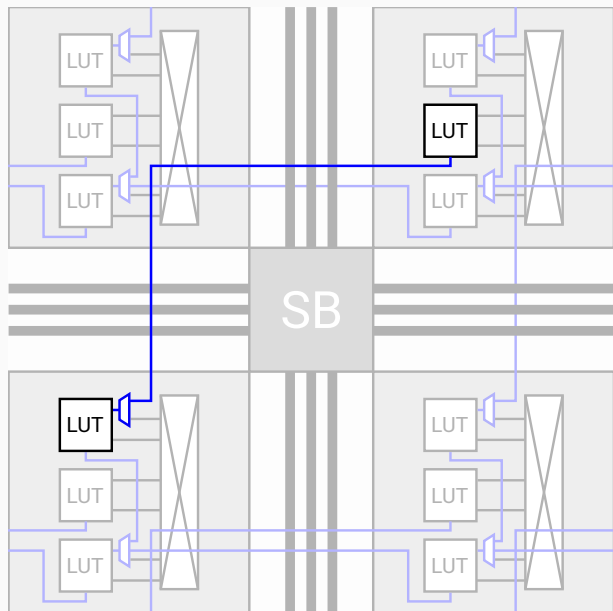
---



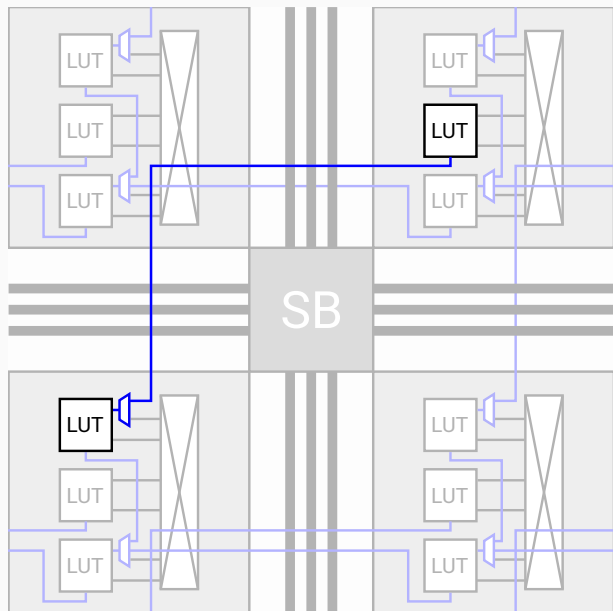
- LUT-to-LUT connections



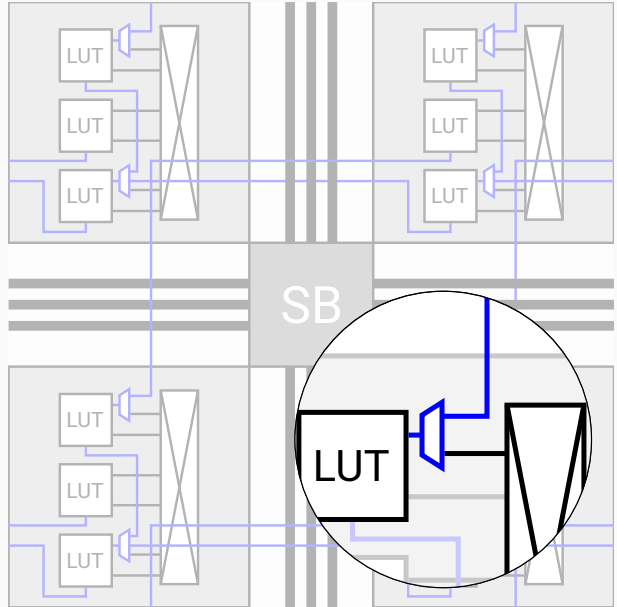
- LUT-to-LUT connections
- Can span multiple clusters



- LUT-to-LUT connections
- Can span multiple clusters
- Optionally used  $\implies$  keeps all the flexibility of the programmable interconnect



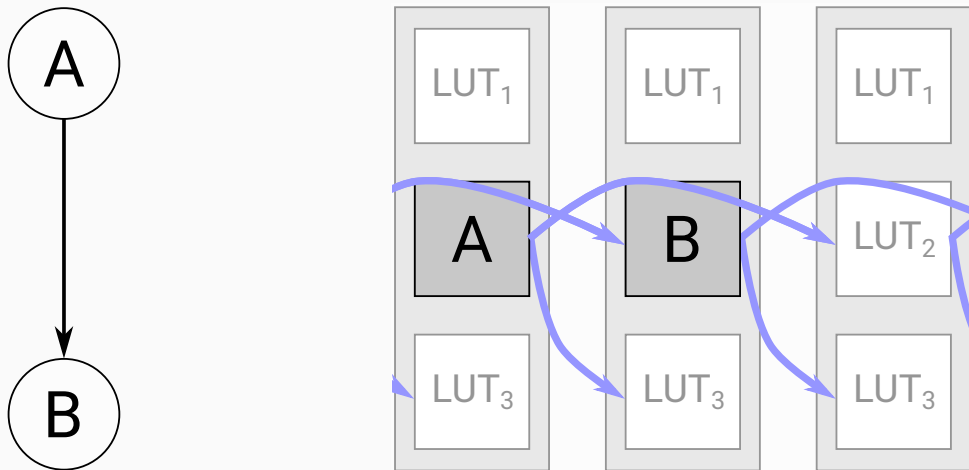
- LUT-to-LUT connections
- Can span multiple clusters
- Optionally used  $\implies$  keeps all the flexibility of the programmable interconnect



# Motivation

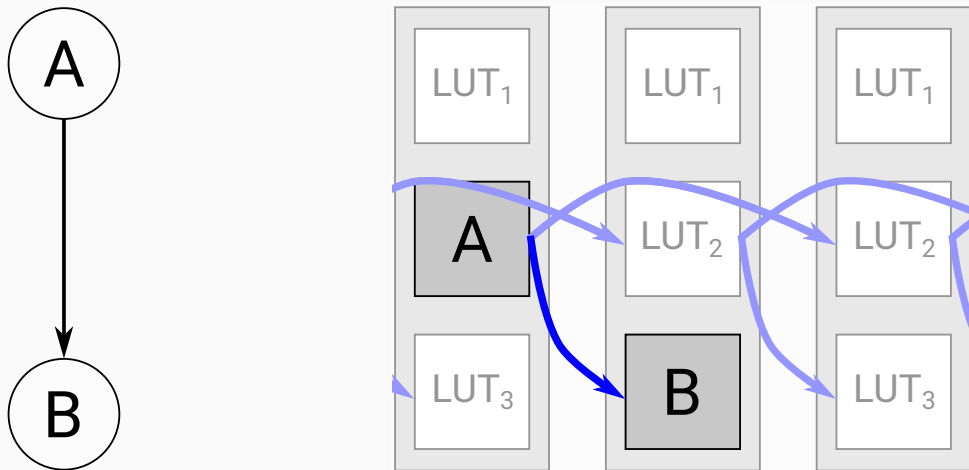
---

# FPGA'20: Swapping LUTs within Clusters



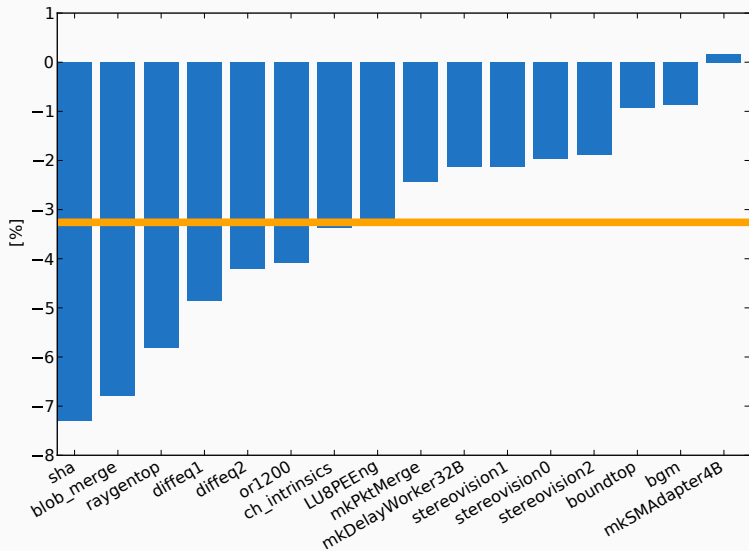
No direct connection between A and B

# FPGA'20: Swapping LUTs within Clusters



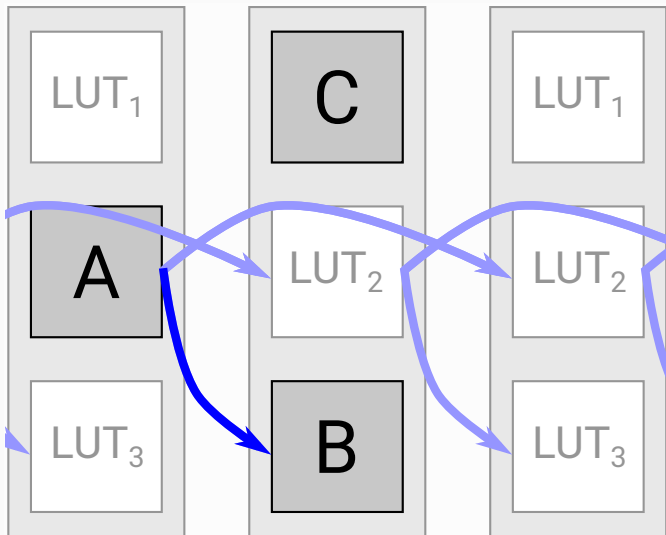
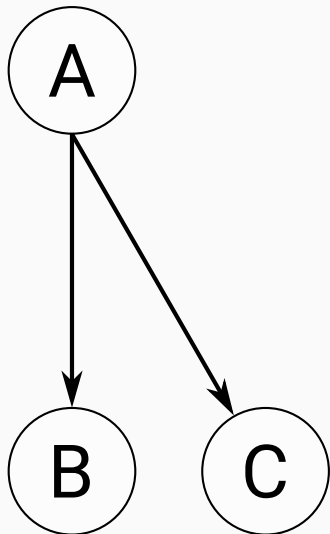
Direct connection between A and B

# FPGA'20: Delay Improvement due to Direct Connections

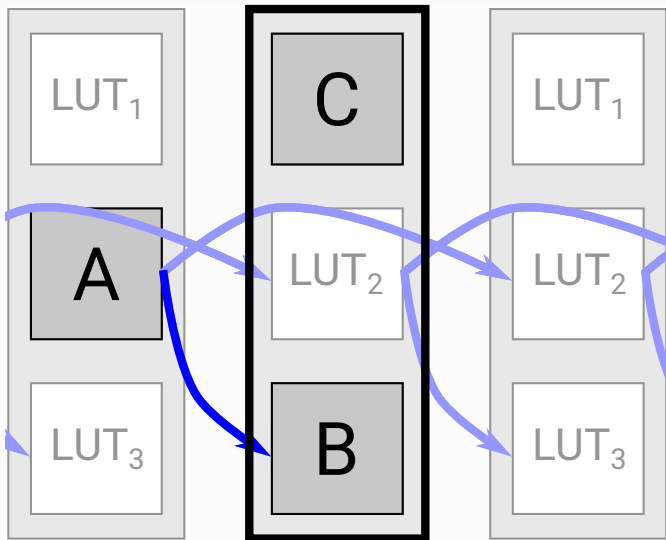
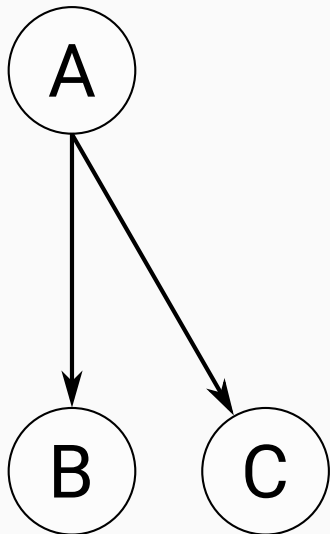




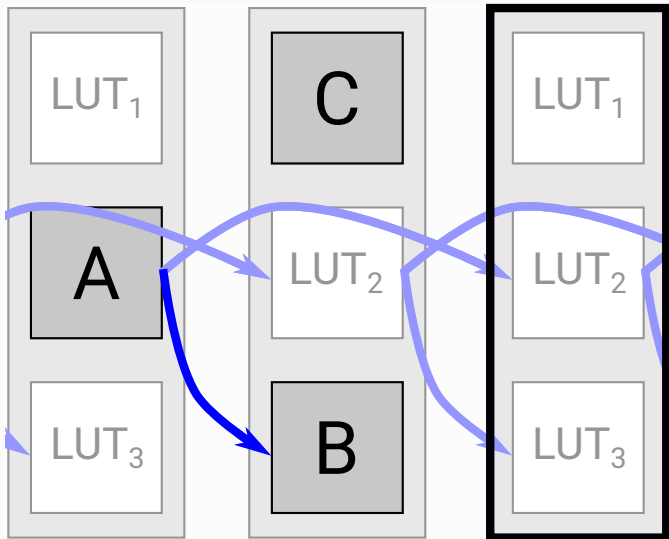
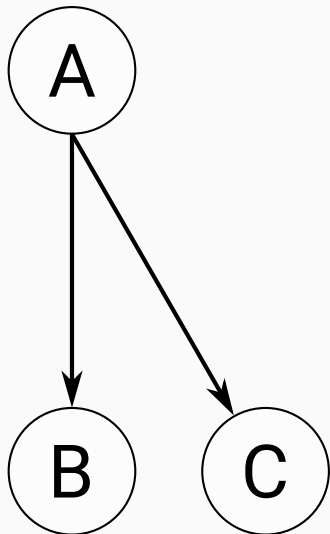
# FPGA'20: Missed Opportunities



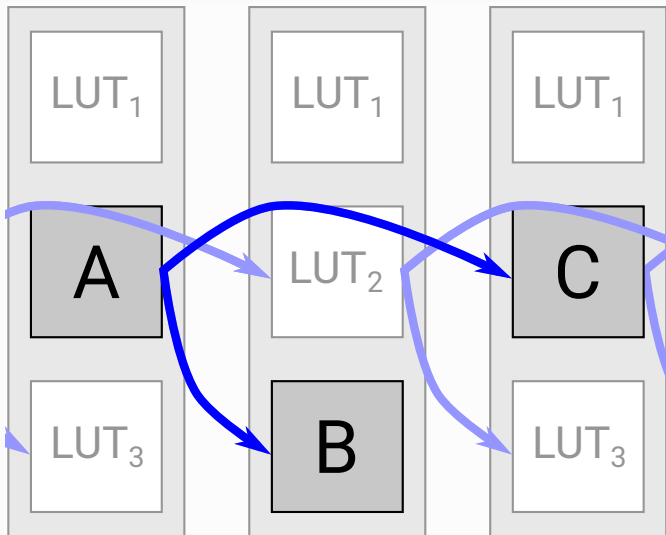
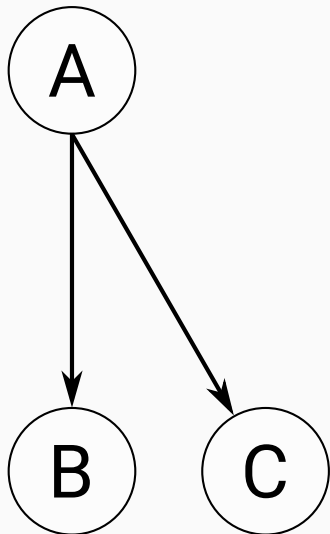
# FPGA'20: Missed Opportunities



# FPGA'20: Missed Opportunities



# FPGA'20: Missed Opportunities



How much could we gain?

# FPGA'20: Missed Opportunities

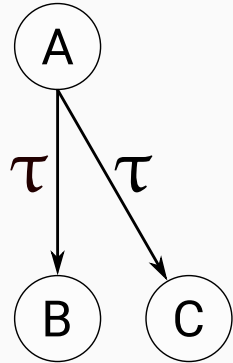
$$\tau = \langle t_d(u, v) \rangle,$$

$\forall(u, v) : (u, v)$  is a direct connection

# FPGA'20: Missed Opportunities

$$\tau = \langle t_d(u, v) \rangle,$$

$\forall(u, v) : (u, v)$  is a direct connection

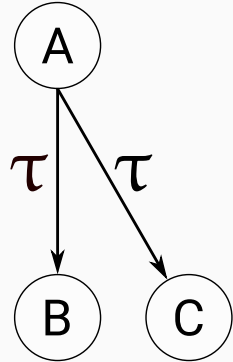


# FPGA'20: Missed Opportunities

$$\tau = \langle t_d(u, v) \rangle,$$

$\forall(u, v) : (u, v)$  is a direct connection

~ 19% lower geomean delay



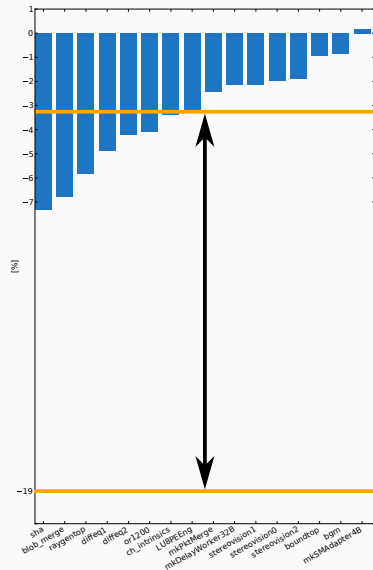


Unlikely to meet in practice...

# FPGA'20: Missed Opportunities

Unlikely to meet in practice...

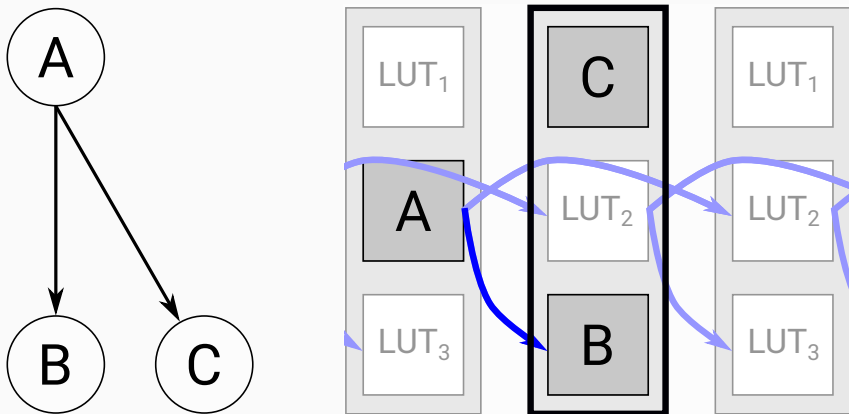
But, leaves a big margin for improvement



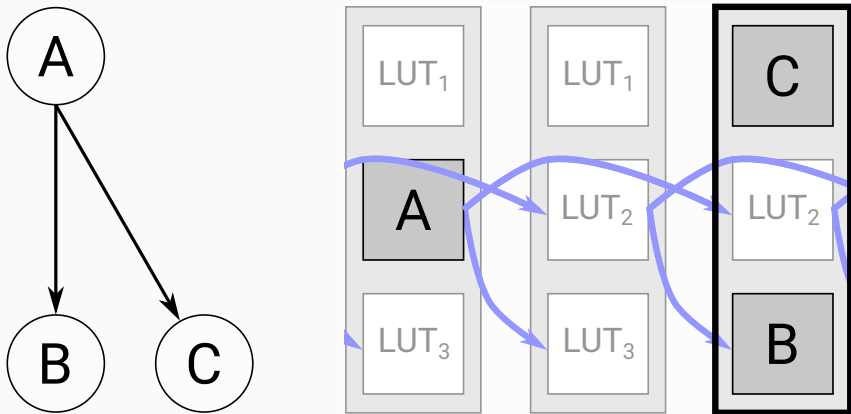
# General Approach

---

# Placing Clusters is not Sufficient



# Placing Clusters is not Sufficient

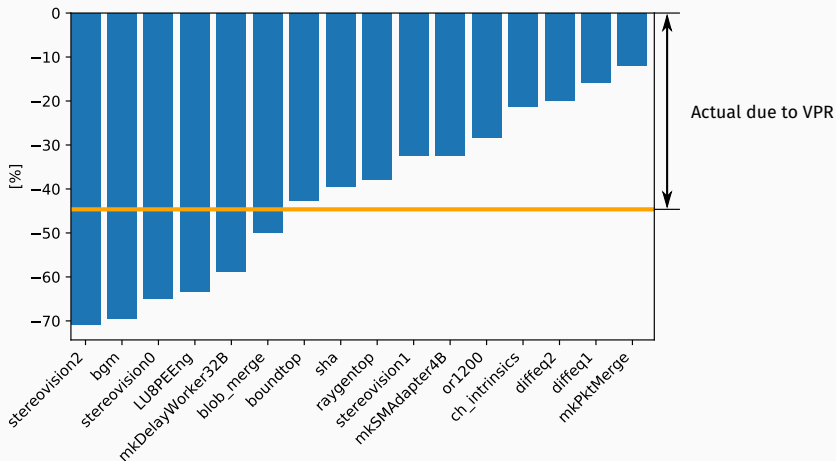


# Flat placement of LUTs

# Flat placement of LUTs

An order of magnitude more  
placeable objects and placement positions

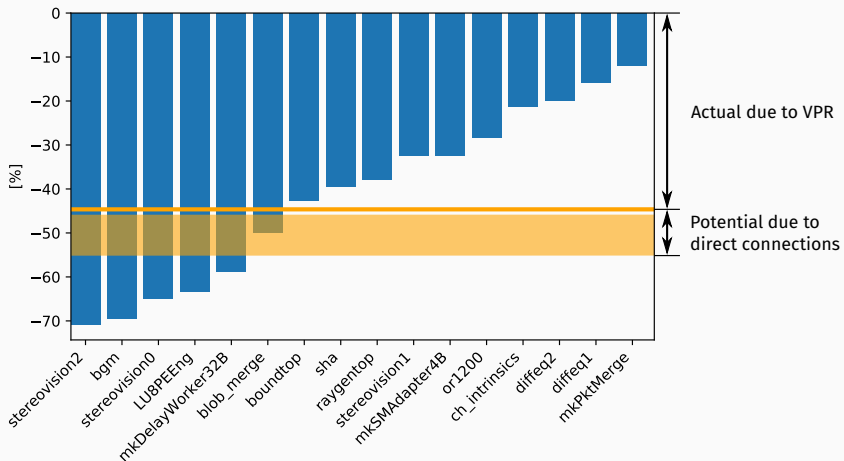
# Generic vs Dedicated Placement



Delay improvement over initial random placement



# Generic vs Dedicated Placement

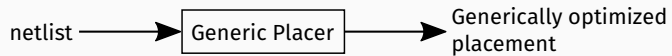


Delay improvement over initial random placement

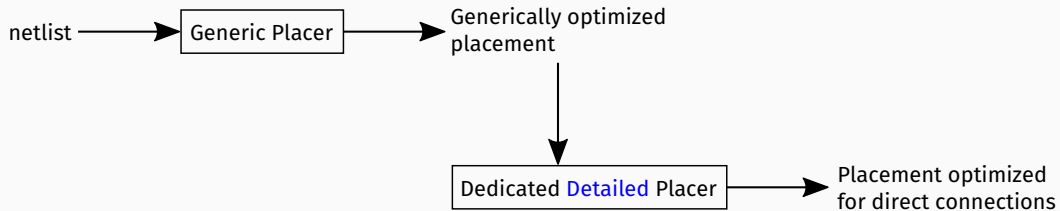
## When to Consider Direct Connections?



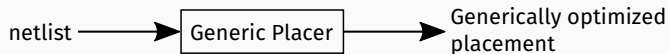
# When to Consider Direct Connections?



# When to Consider Direct Connections?



# When to Consider Direct Connections?



Dedicated **Detailed** Placer

nized  
ctions

# Placement Algorithm

---

# Timing-Driven Detailed Placement

0. All nodes (LUTs) are assigned a starting position

# Timing-Driven Detailed Placement

0. All nodes (LUTs) are assigned a starting position
1. Select a subset of nodes



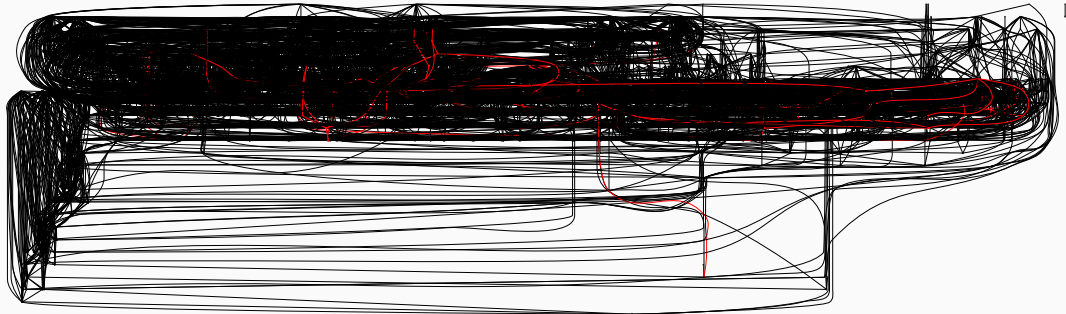
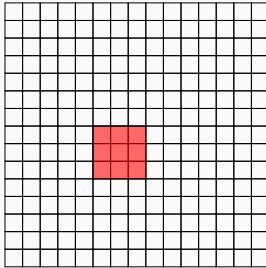
# Timing-Driven Detailed Placement

0. All nodes (LUTs) are assigned a starting position
1. Select a subset of nodes
2. Move them to reduce the critical path delay

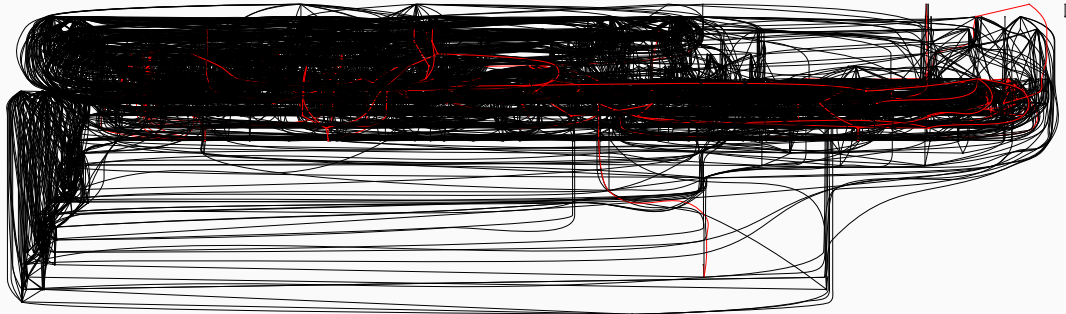
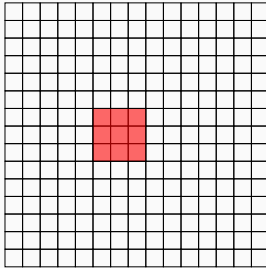
Which Nodes to Move?

---

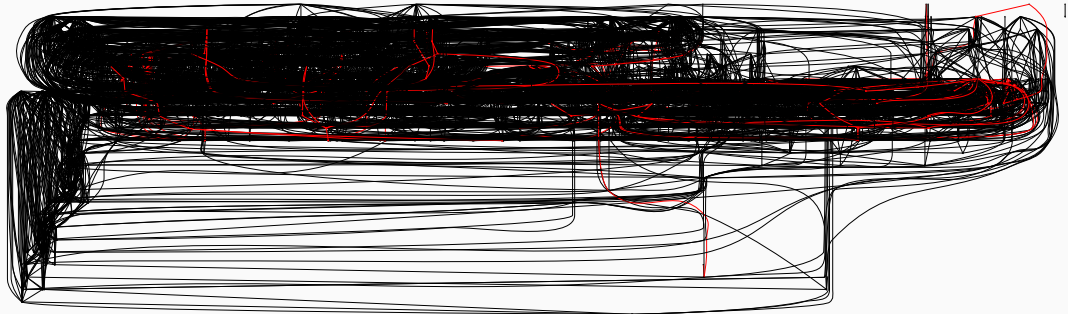
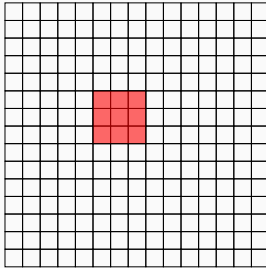
# Determining Movable Nodes: Sliding Window



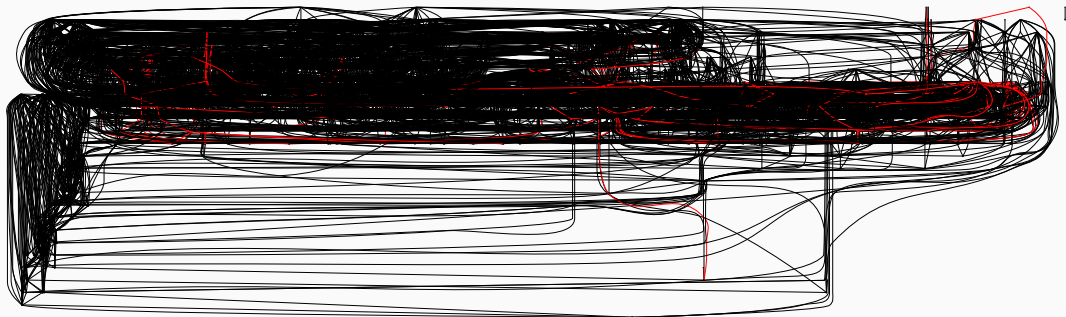
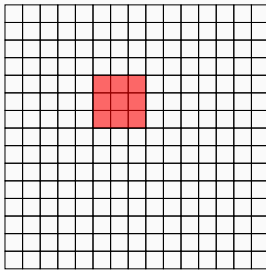
# Determining Movable Nodes: Sliding Window



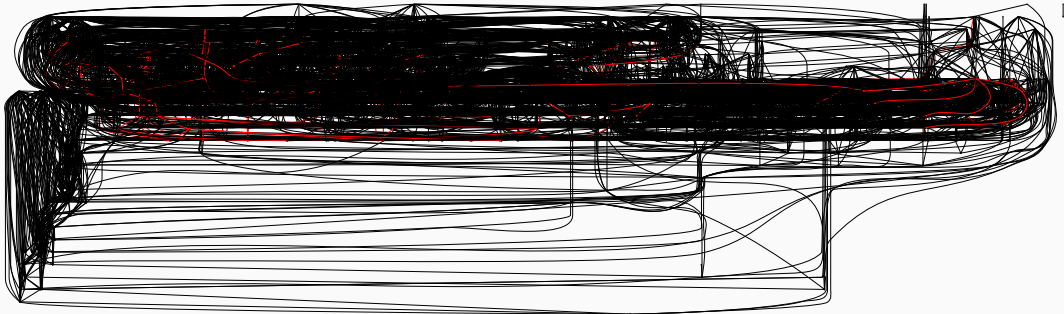
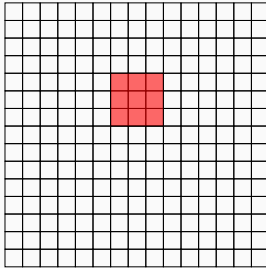
# Determining Movable Nodes: Sliding Window



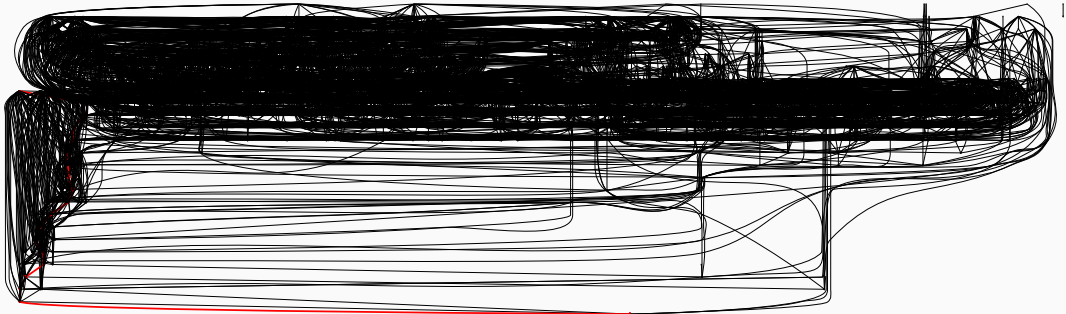
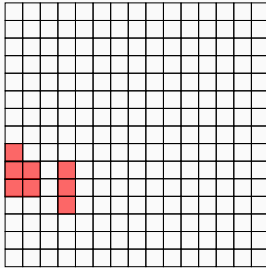
# Determining Movable Nodes: Sliding Window



# Determining Movable Nodes: Sliding Window



# Determining Movable Nodes: Critical Path





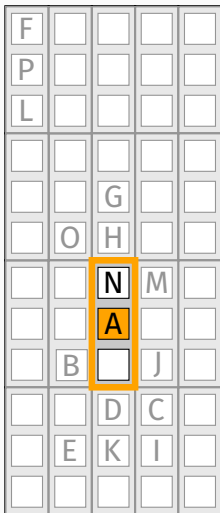
# Generalization

---

# Movement Constraints

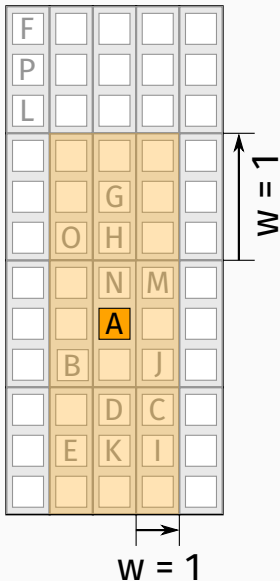
F				
P				
L				
		G		
	O	H		
		N	M	
		A		
	B		J	
		D	C	
	E	K	I	

# Movement Constraints



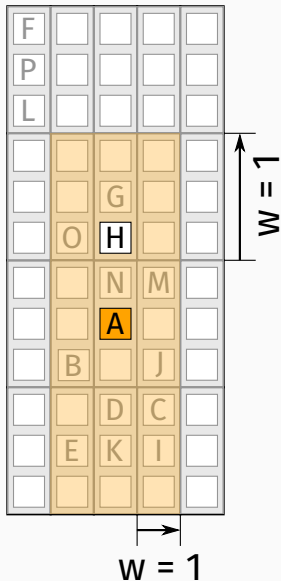
- Each node can move to any position in the  $w$ -bounded square around its starting cluster

# Movement Constraints



- Each node can move to any position in the  $w$ -bounded square around its starting cluster

# Movement Constraints



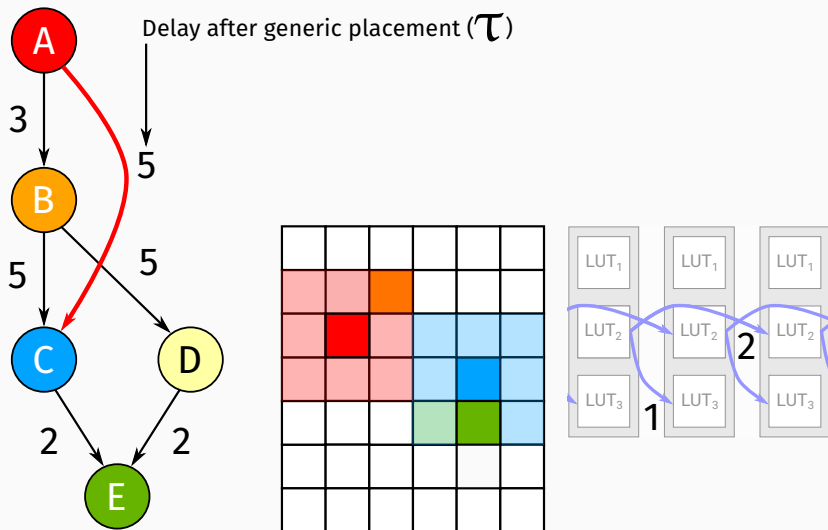
- Each node can move to any position in the  $w$ -bounded square around its starting cluster
- Overlaps with stationary nodes removed by postprocessing

## Improving Connection Delays

Each circuit connection  $(u, v)$  has initial delay  $\tau_{u,v}$

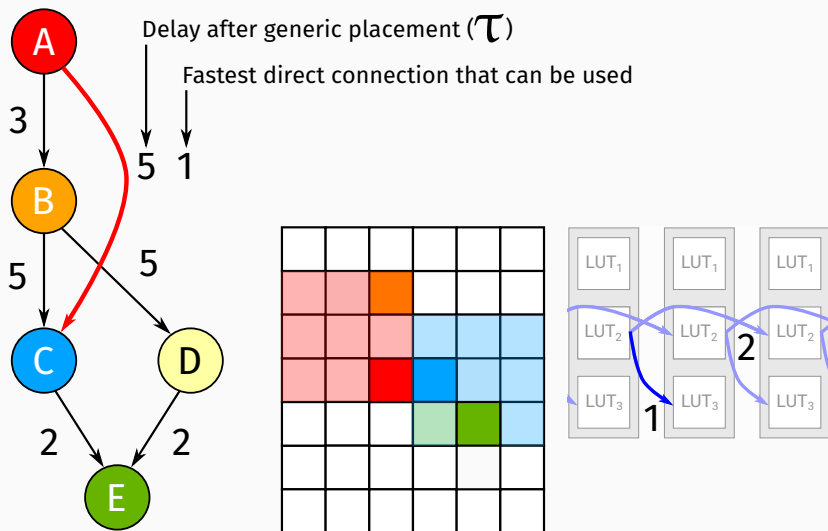
Implementing by a direct connection can improve it  
by  $0 \leq \mathit{imp}_{u,v} \leq l_{u,v} = \mathit{const}$ .

# Improving Connection Delays

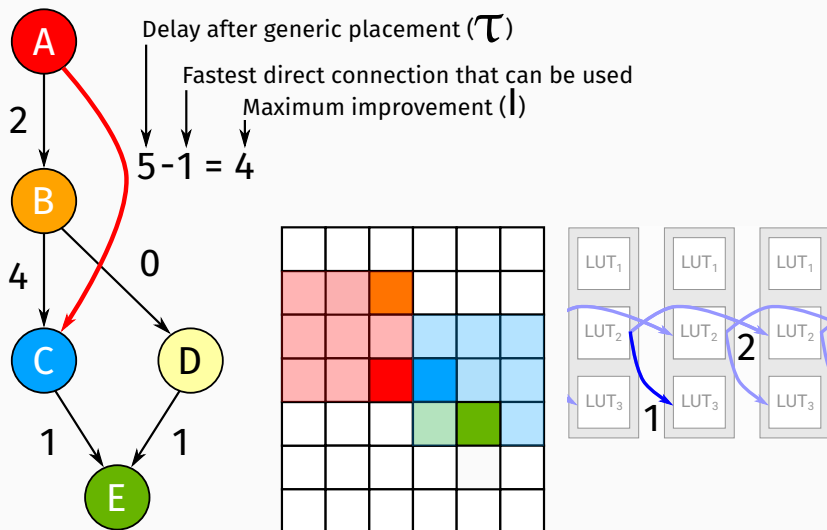




# Improving Connection Delays



# Improving Connection Delays



# Improving Connection Delays

1. Assign *imp*-variables values, s.t. critical path delay  $\leq$  some target  $D$

# Improving Connection Delays

1. Assign *imp*-variables values, s.t. critical path delay  $\leq$  some target  $D$
2. All nodes incident to an edge with  $imp \neq 0$  are movable

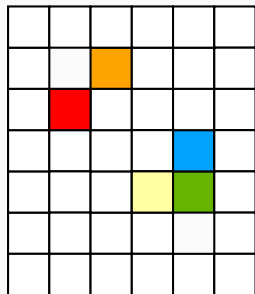
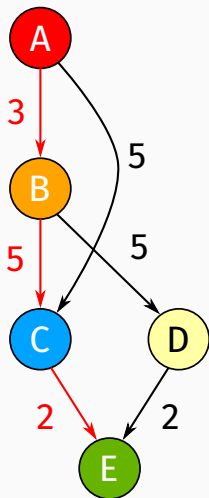
# Improving Connection Delays

1. Assign *imp*-variables values, s.t. critical path delay  $\leq$  some target  $D$
2. All nodes incident to an edge with  $imp \neq 0$  are movable

$$\min |\{(u, v) : imp_{u,v} \neq 0\}|$$

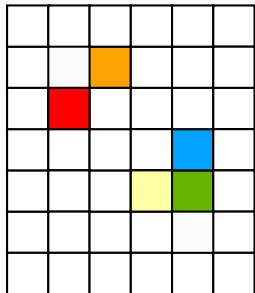
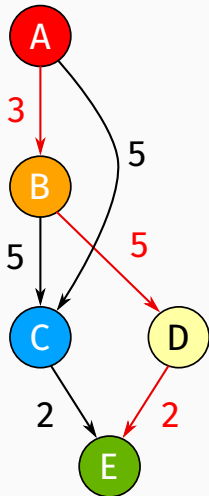
# Improving Connection Delays: An Example

Two critical paths with delay 10



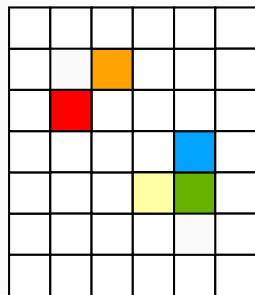
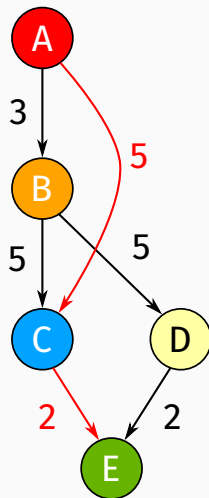
# Improving Connection Delays: An Example

Two critical paths with delay 10



# Improving Connection Delays: An Example

Two critical paths with delay 10  
One path with delay 7



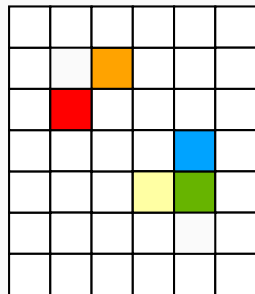
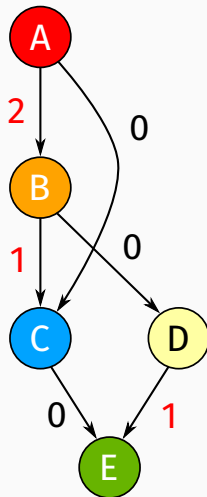


# Improving Connection Delays: An Example

Two critical paths with delay 10

One path with delay 7

$D = 7$



# Improving Connection Delays: Selection LP

1. ... s.t. critical path delay  $\leq D$
2. ...

# Improving Connection Delays: Selection LP

1. ... s.t. critical path delay  $\leq D$

2. ...

$$\text{s.t. } t_{u,v} = \tau_{u,v} - \text{imp}_{u,v}$$

$$ta_v \geq ta_u + t_{u,v}$$

$$ta_u \leq ta_{max}$$

$$ta_{max} \leq D$$

# Improving Connection Delays: Selection LP

1. ... s.t. critical path delay  $\leq D$

2. ...

$$\min \{(u, v) : imp_{u,v} \neq 0\}$$

$$\text{s.t. } t_{u,v} = \tau_{u,v} - imp_{u,v}$$

$$ta_v \geq ta_u + t_{u,v}$$

$$ta_u \leq ta_{max}$$

$$ta_{max} \leq D$$

# Improving Connection Delays: Selection LP

1. ... s.t. critical path delay  $\leq D$

2. ...

$$\min |\{(u, v) : imp_{u,v} \neq 0\}|$$

$$\min \sum_{(u,v)} imp_{u,v}$$

$$\text{s.t. } t_{u,v} = \tau_{u,v} - imp_{u,v}$$

$$ta_v \geq ta_u + t_{u,v}$$

$$ta_u \leq ta_{max}$$

$$ta_{max} \leq D$$

# Improving Connection Delays: Selection LP

$$\min \sum_{(u,v)} imp_{u,v}$$

$$\text{s.t. } t_{u,v} = \tau_{u,v} - imp_{u,v}$$

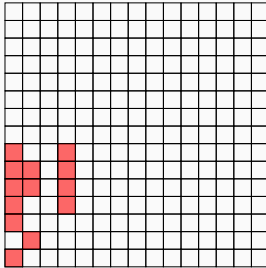
$$ta_v \geq ta_u + t_{u,v}$$

$$ta_u \leq ta_{max}$$

$$ta_{max} \leq D$$

<sup>1</sup>Hambrusch and Tu, "Edge weight reduction problems in directed acyclic graphs", *J. Algorithms*, 1997

# Determining Movable Nodes: Selection LP



## How to Move the Selected Nodes?

---



Heuristic Methods:

Heuristic Methods:

...

...

...

Heuristic Methods:

...

...

...

Exact Methods:

Heuristic Methods:

...

...

...

Exact Methods:

SAT

SMT

ILP

Exact Methods:

SAT

SMT

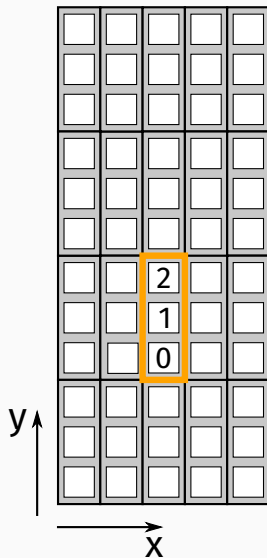
ILP

Exact Methods:

ILP

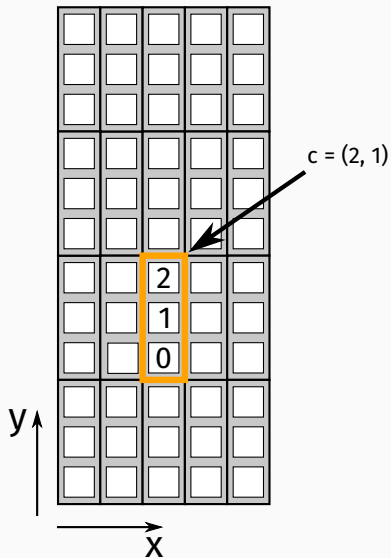
# Naive ILP: Describing a Particular Placement

Cluster position:  $c = (x, y)$



# Naive ILP: Describing a Particular Placement

Cluster position:  $c = (x, y)$

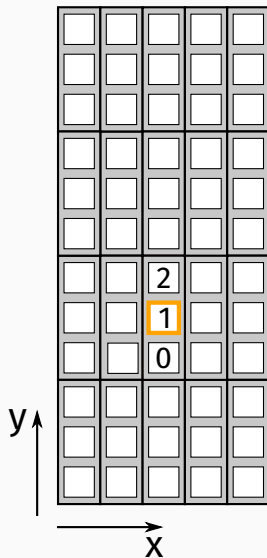




# Naive ILP: Describing a Particular Placement

Cluster position:  $c = (x, y)$

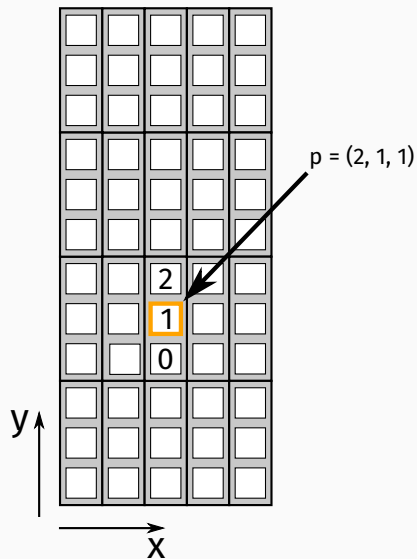
LUT position:  $p = (x, y, i)$



# Naive ILP: Describing a Particular Placement

Cluster position:  $c = (x, y)$

LUT position:  $p = (x, y, i)$

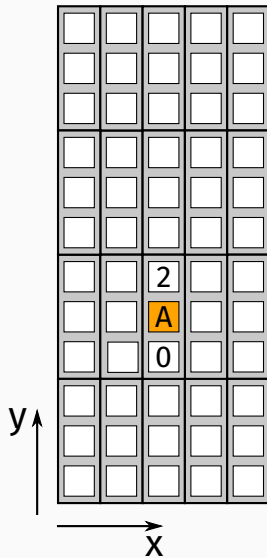


# Naive ILP: Describing a Particular Placement

Cluster position:  $c = (x, y)$

LUT position:  $p = (x, y, i)$

Introduce:  $x_{u,p} \in \{0, 1\}$

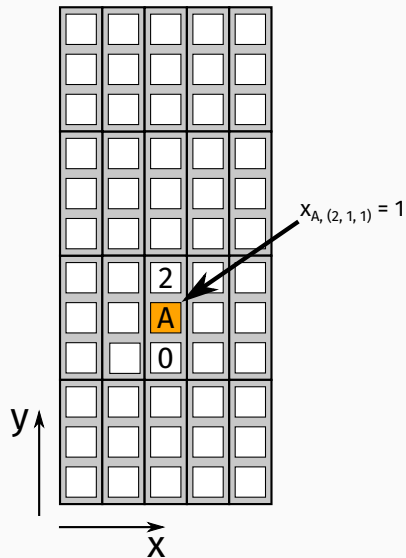


# Naive ILP: Describing a Particular Placement

Cluster position:  $c = (x, y)$

LUT position:  $p = (x, y, i)$

Introduce:  $x_{u,p} \in \{0, 1\}$

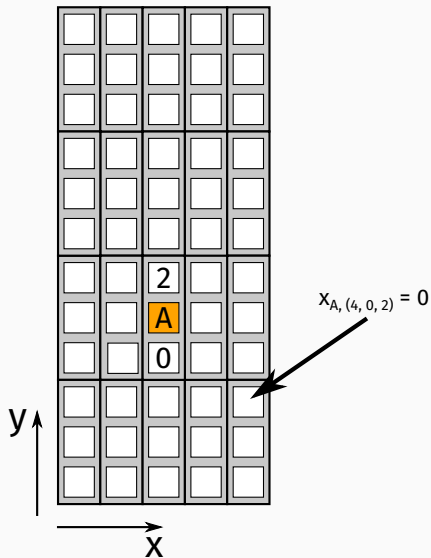


# Naive ILP: Describing a Particular Placement

Cluster position:  $c = (x, y)$

LUT position:  $p = (x, y, i)$

Introduce:  $x_{u,p} \in \{0, 1\}$

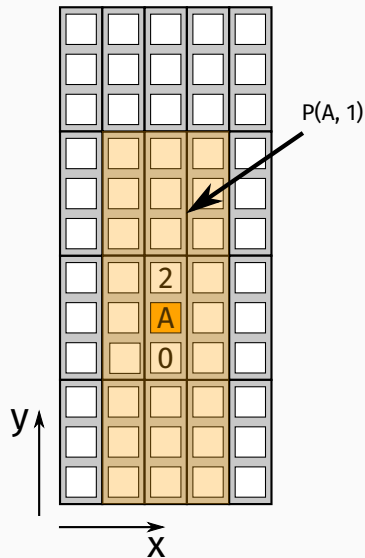


# Naive ILP: Describing a Particular Placement

Cluster position:  $c = (x, y)$

LUT position:  $p = (x, y, i)$

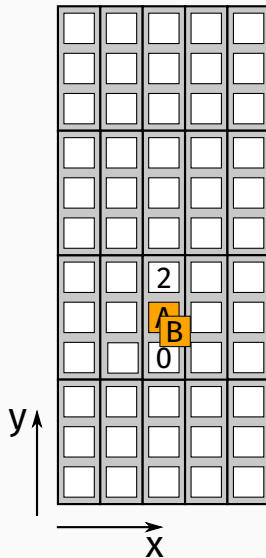
Introduce:  $x_{u,p} \in \{0, 1\}, \forall p \in P(u, w)$



# Naive ILP: Describing Any Legal Placement

No overlaps between movable nodes:

$$\sum_{u \in V_m} x_{u,p} \leq 1, \forall p$$



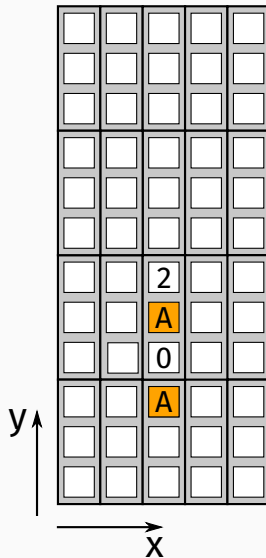
# Naive ILP: Describing Any Legal Placement

No overlaps between movable nodes:

$$\sum_{u \in V_m} x_{u,p} \leq 1, \forall p$$

Each node uniquely placed:

$$\sum_{p \in P(u,w)} x_{u,p} = 1, \forall u$$





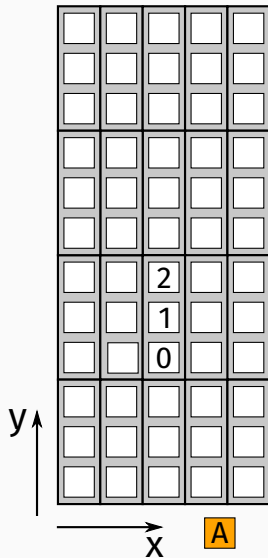
# Naive ILP: Describing Any Legal Placement

No overlaps between movable nodes:

$$\sum_{u \in V_m} x_{u,p} \leq 1, \forall p$$

Each node uniquely placed:

$$\sum_{p \in P(u,w)} x_{u,p} = 1, \forall u$$



# Naive ILP: Timing

Arrival times: same as Selection LP

# Naive ILP: Timing

Arrival times: same as Selection LP

$$\text{Connection delay: } t_{u,v} = \sum_{p_u \in P(u,w), p_v \in P(v,w)} \tau_{p_u, p_v} X_{u, p_u} X_{v, p_v}$$

# Naive ILP: Timing

Arrival times: same as Selection LP

Connection delay:  $t_{u,v} = \sum_{p_u \in P(u,w), p_v \in P(v,w)} \tau_{p_u, p_v} X_{u, p_u} X_{v, p_v}$

## Naive ILP: Encoding Efficiency

$$t_{u,v} = \sum_{p_u \in P(u,w), p_v \in P(v,w)} \tau_{p_u, p_v} x_{u, p_u} x_{v, p_v}$$

## Naive ILP: Encoding Efficiency

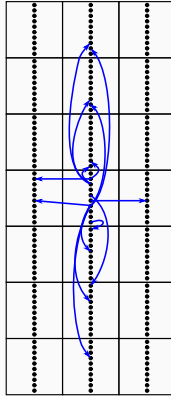
$$t_{u,v} = \sum_{p_u \in \underbrace{P(u,w), p_v \in P(v,w)}_{(2w+1)^2 N}} \tau_{p_u, p_v} x_{u, p_u} x_{v, p_v}$$

## Naive ILP: Encoding Efficiency

$$t_{u,v} = \sum_{\substack{p_u \in P(u,w), p_v \in P(v,w)}} \tau_{p_u, p_v} \underbrace{X_{u,p_u} X_{v,p_v}}_{((2w+1)^2 N)^2}$$

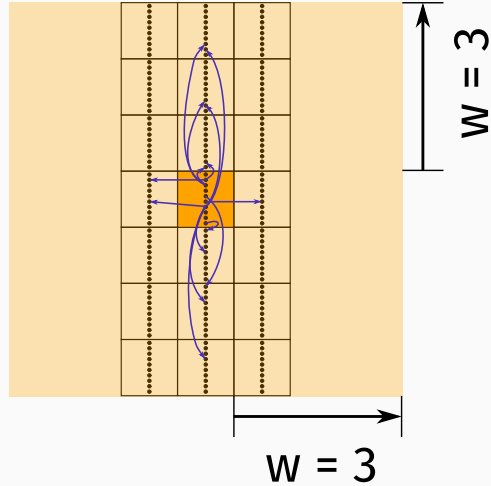
$(2w+1)^2 N$                        $((2w+1)^2 N)^2$

# An Example Target Architecture (FPGA'20)





# An Example Target Architecture (FPGA'20)

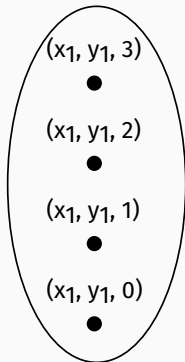


## Naive ILP: Encoding Efficiency

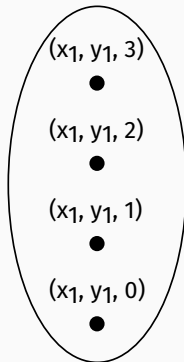
$$t_{u,v} = \sum_{\substack{p_u \in P(u,w), p_v \in P(v,w)}} \tau_{p_u, p_v} \underbrace{X_{u,p_u} X_{v,p_v}}_{240,100}$$

490

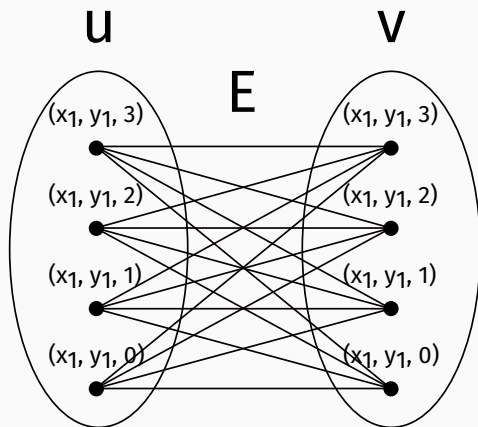
**U**



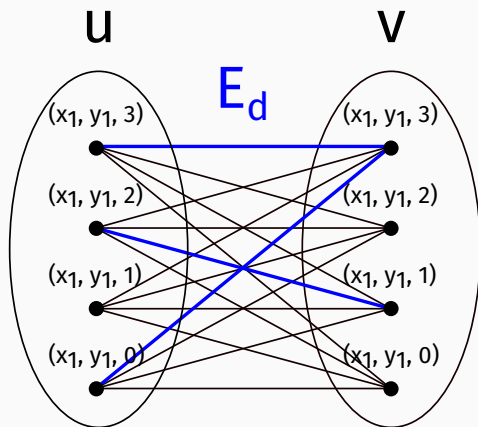
**V**



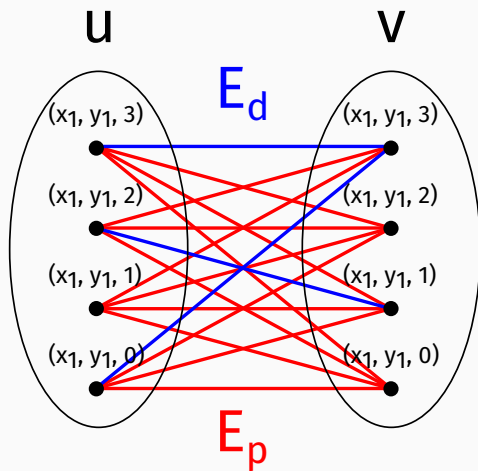
# Improved ILP



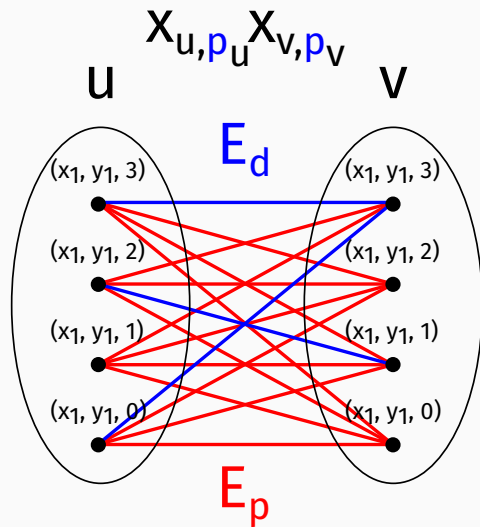
# Improved ILP



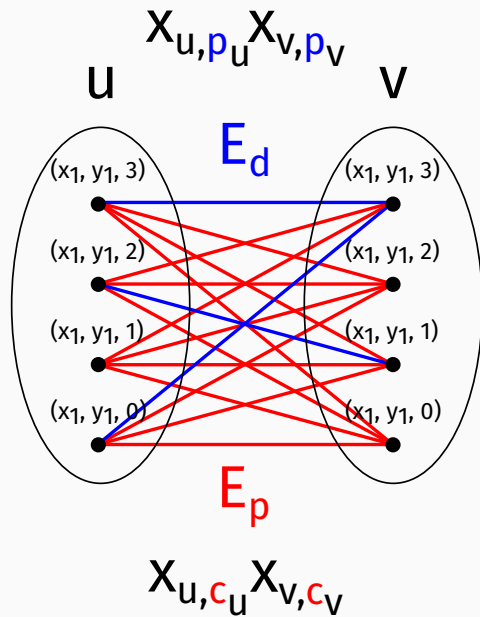
# Improved ILP



# Improved ILP

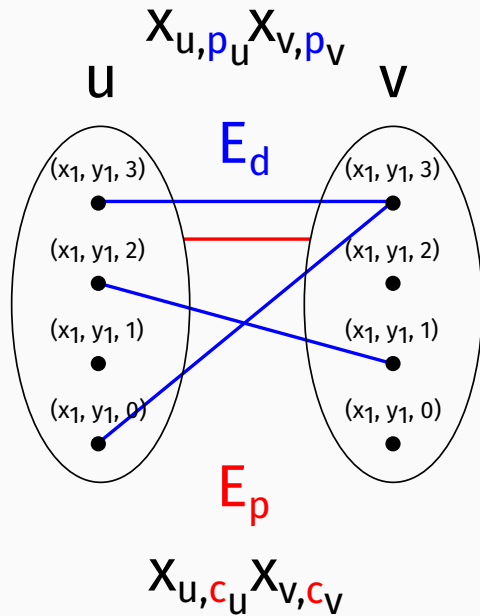


# Improved ILP





# Improved ILP



$$t_{u,v} = \sum_E \tau_{p_u, p_v} x_{u, p_u} x_{v, p_v}$$

$$t_{u,v} = \sum_E \tau_{p_u, p_v} x_{u, p_u} x_{v, p_v}$$

$$t_{u,v} = \sum_{E_d} \tau_{p_u, p_v} x_{u, p_u} x_{v, p_v}$$

$$t_{u,v} = \sum_E \tau_{p_u, p_v} x_{u, p_u} x_{v, p_v}$$

$$t_{u,v} = \sum_{E_d} \tau_{p_u, p_v} x_{u, p_u} x_{v, p_v} + \sum_{E_p} \tau_{c_u, c_v} x_{u, c_u} x_{v, c_v}$$

$$t_{u,v} = \sum_E \tau_{p_u,p_v} x_{u,p_u} x_{v,p_v}$$

$$t_{u,v} = y \sum_{E_d} \tau_{p_u,p_v} x_{u,p_u} x_{v,p_v} + (1 - y) \sum_{E_p} \tau_{c_u,c_v} x_{u,c_u} x_{v,c_v}$$

# Improved ILP

$$t_{u,v} = \sum_E \tau_{p_u,p_v} x_{u,p_u} x_{v,p_v}$$

$$t_{u,v} = y \sum_{E_d} \tau_{p_u,p_v} x_{u,p_u} x_{v,p_v} + (1 - y) \sum_{E_p} \tau_{c_u,c_v} x_{u,c_u} x_{v,c_v}$$

$\uparrow$   
 $((2w + 1)^2 N)^2$

# Improved ILP

$$t_{u,v} = \sum_E \tau_{p_u, p_v} x_{u, p_u} x_{v, p_v}$$

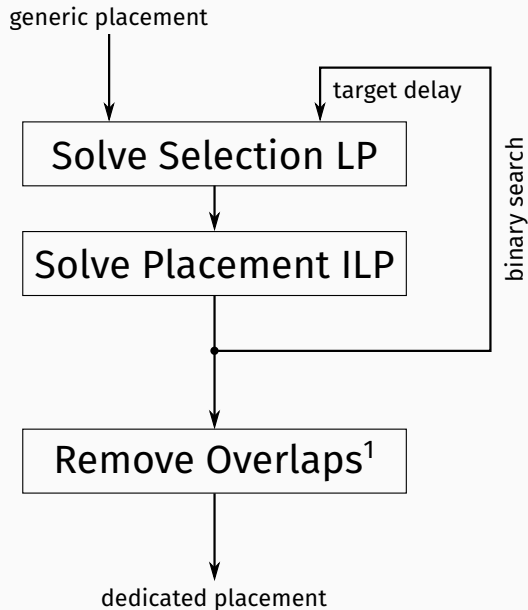
$$t_{u,v} = y \sum_{E_d} \tau_{p_u, p_v} x_{u, p_u} x_{v, p_v} + (1 - y) \sum_{E_p} \tau_{c_u, c_v} x_{u, c_u} x_{v, c_v}$$

$\uparrow$   
 $((2w + 1)^2 M)^2$

## Complete Flow

---





<sup>1</sup>Darav et al., "Multi-commodity flow-based spreading in a commercial analytic placer", *FPGA'19*

# Experimental Setup

---

Almost the same as FPGA'20

Almost the same as FPGA'20

- Architecture: best found in FPGA'20
  - 14 direct connections, all crossing clusters
  - 10 6-LUT cluster
  - 40 inputs
  - Complete crossbar
  - No carry chains

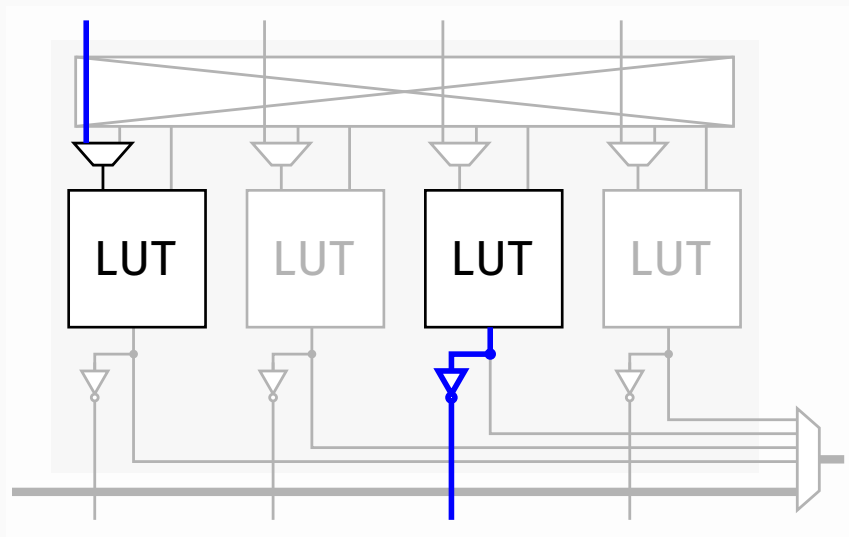
Almost the same as FPGA'20

- Architecture: best found in FPGA'20
  - 14 direct connections, all crossing clusters
  - 10 6-LUT cluster
  - ~~40~~ 60 inputs
  - Complete crossbar
  - No carry chains

Almost the same as FPGA'20

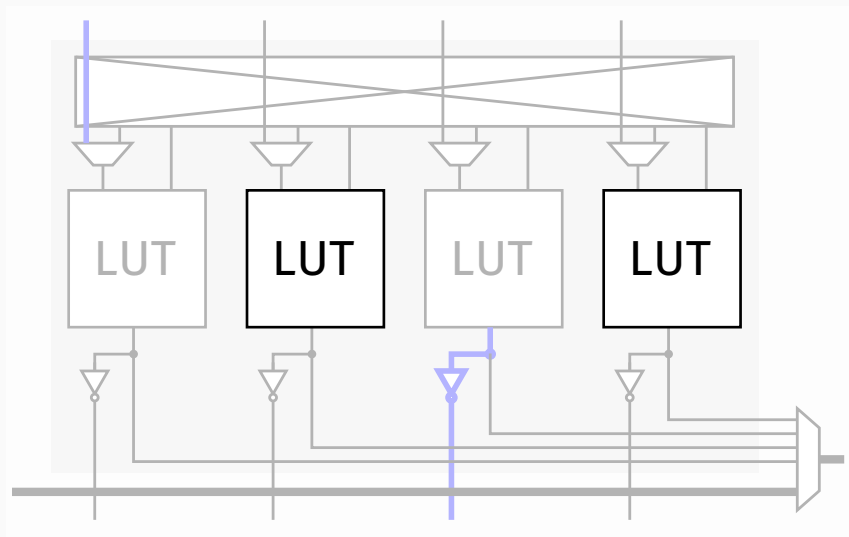
- Architecture: best found in FPGA'20
  - 14 direct connections, all crossing clusters
  - 10 6-LUT cluster
  - ~~40~~ 60 inputs
  - Complete crossbar
  - No carry chains
- VTR 7, with Rubin and DeHon's *delay targeted routing*

## Route-time LUT Permutation



Fixed

## Route-time LUT Permutation



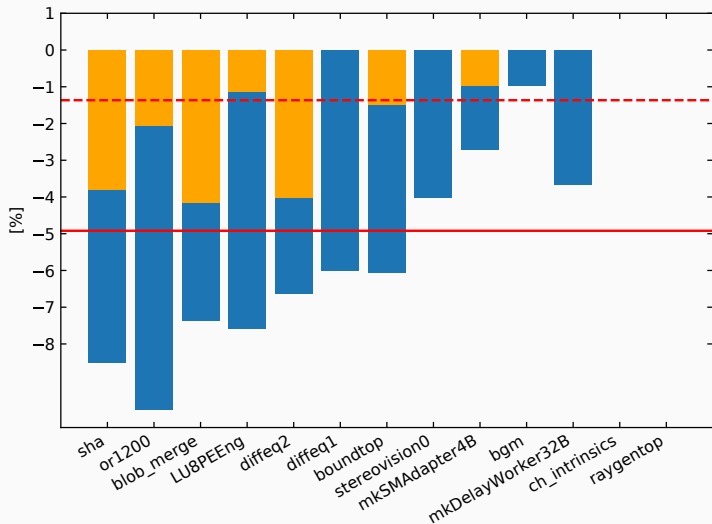
Permutable



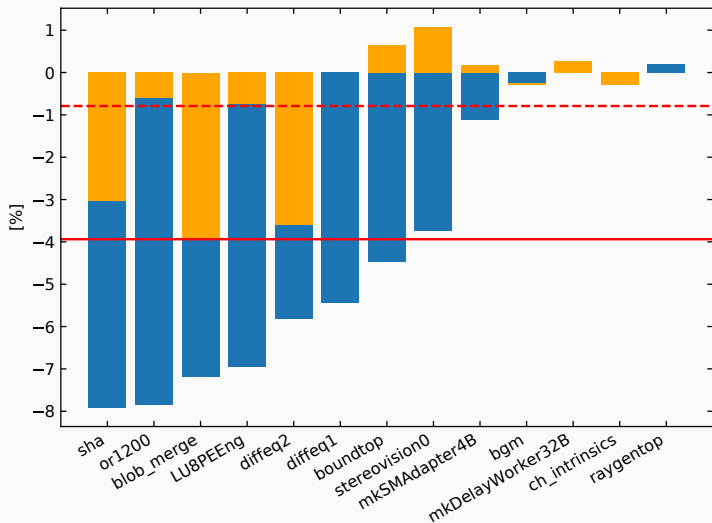
# Results

---

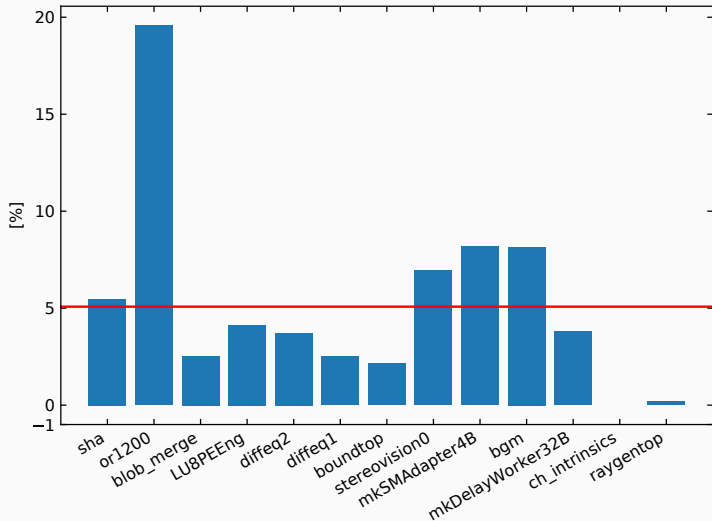
# $w = 1$ vs $w = 0$ Delay Change over Baseline: Postplacement

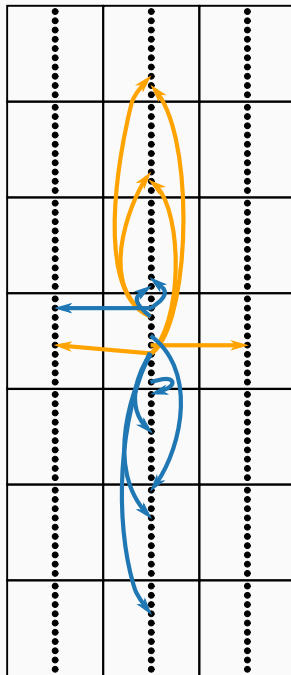
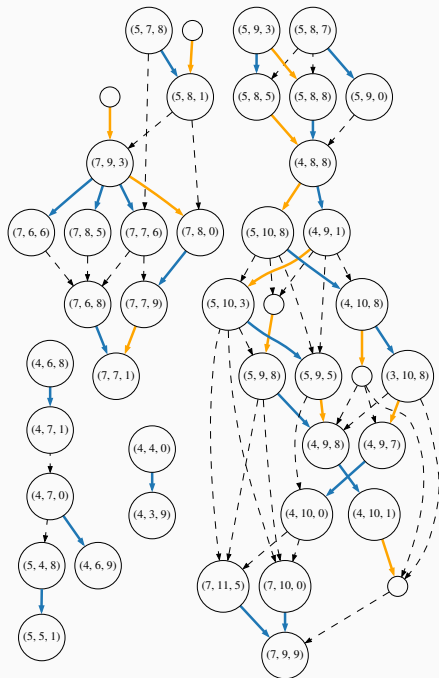


# $w = 1$ vs $w = 0$ Delay Change over Baseline: Postrouting



# $w = 1$ Delay Change over Baseline, all Programmable





We now have an effective dedicated placer for architectures with direct connections

Address scalability issues to extend movement freedom

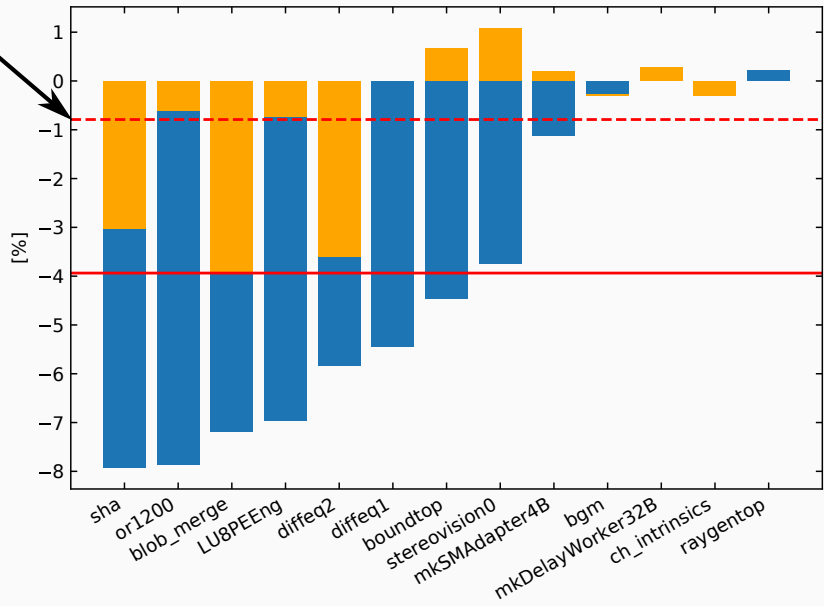
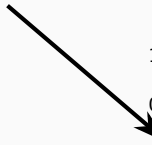
# Future Work

Address scalability issues to extend movement freedom

Or allocate the existing freedom more wisely



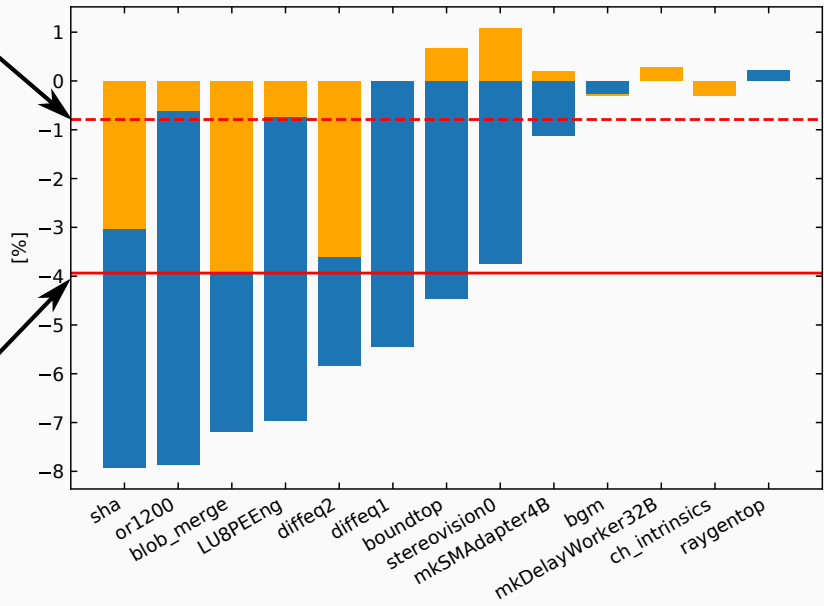
3%?



3%?



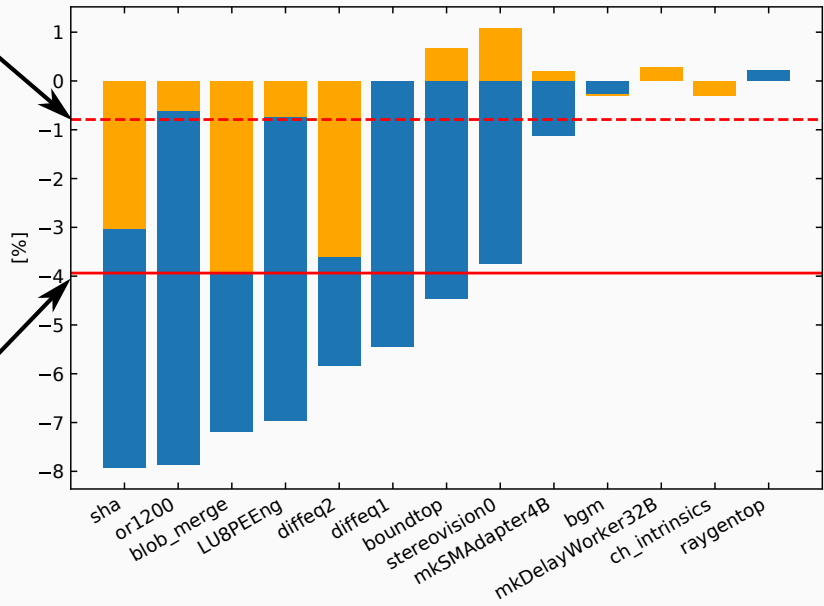
7%?

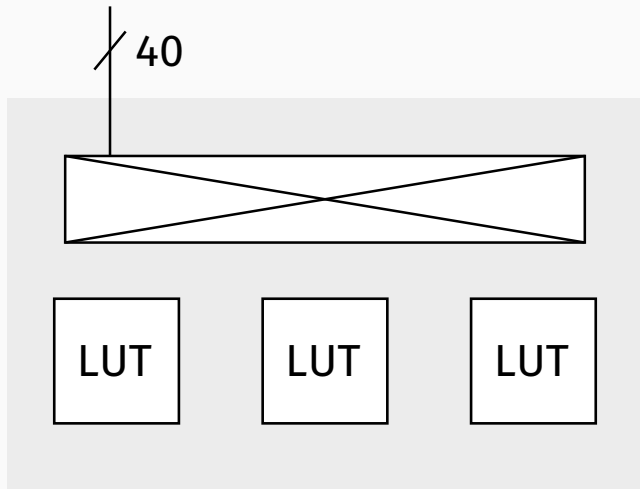


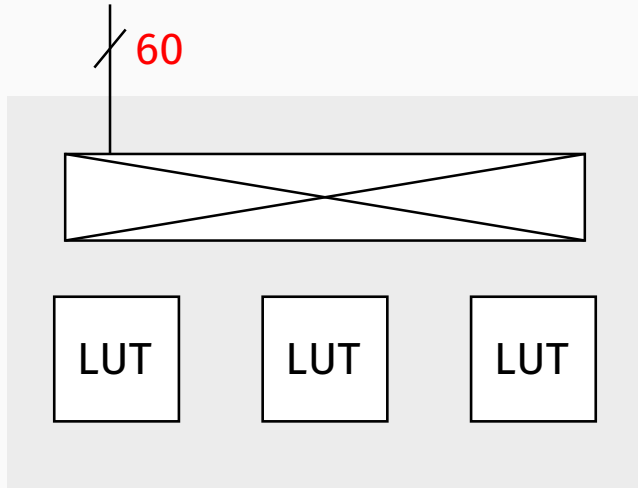
3%?



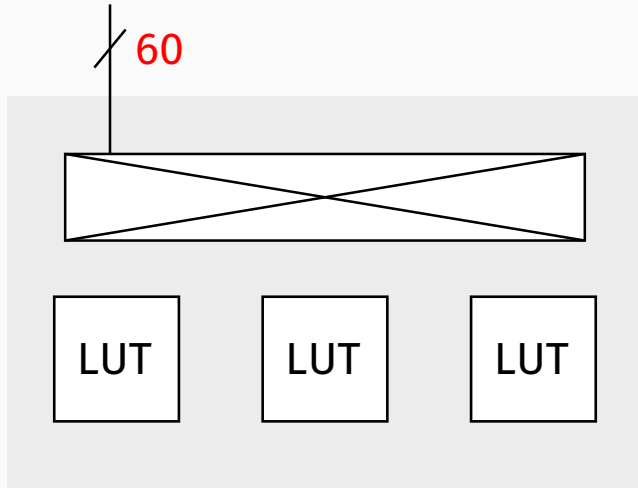
7%?  
12%?



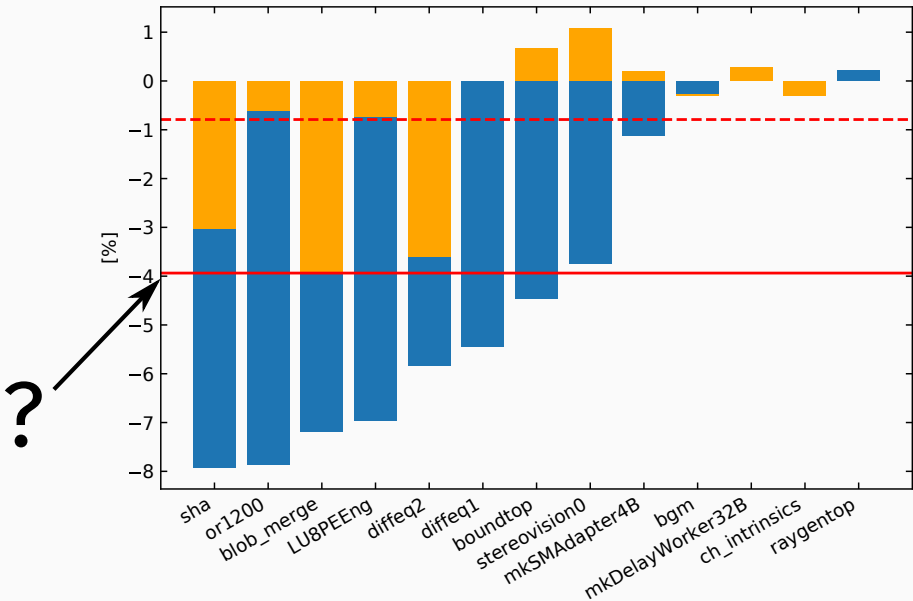




Denser packing



Denser packing + no local direct connections  
⇒ less chance for optimization



# Future Work

Address scalability issues to extend movement freedom

Or allocate the existing freedom more wisely

Extensive architectural exploration



# Thank you for attention

<https://github.com/stefannikolicns/fpl20-placement>