

DREAMPlaceFPGA-PL: An Open-Source GPU-Accelerated Packer-Legalizer for Heterogeneous FPGAs

Rachel Selina Rajarathnam¹, Zixuan Jiang¹, Mahesh A. Iyer², David Z. Pan¹

¹ECE Department, The University of Texas at Austin, TX, USA ²Intel Corporation, CA, USA

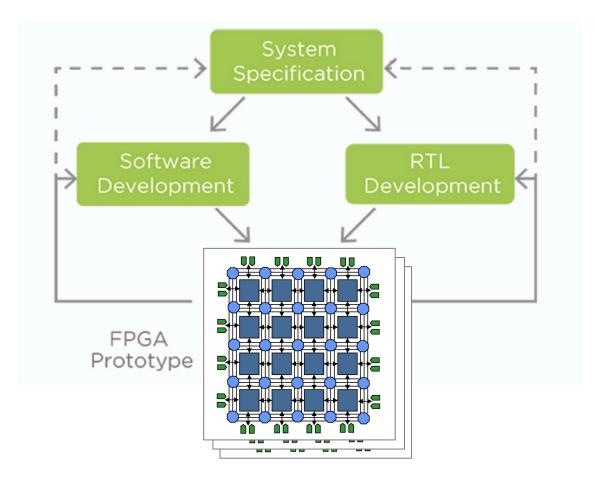




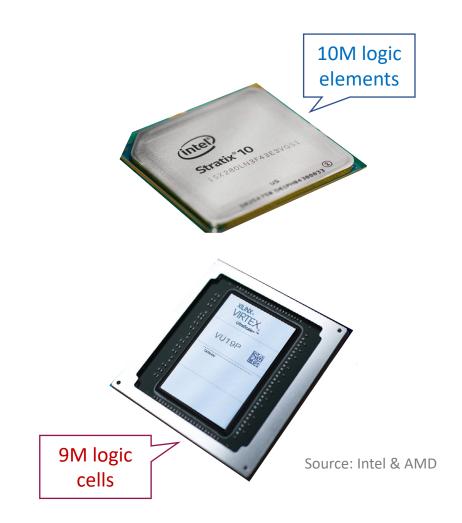
OUTLINE

- Background & Motivation
- Previous Works
- DREAMPlaceFPGA-PL
- Results
- Summary

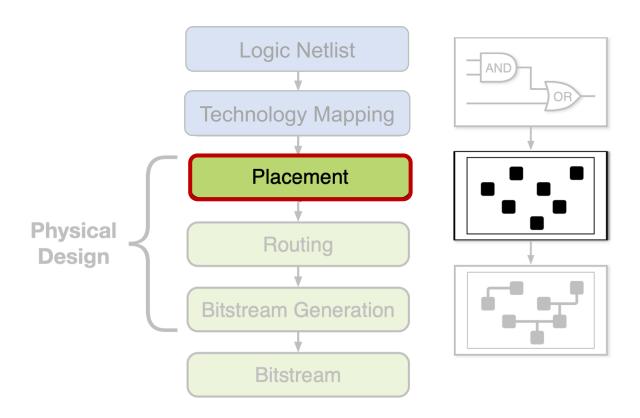
FPGAs IN ASIC PROTOTYPING



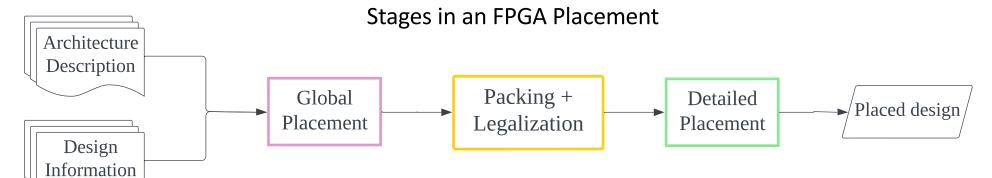
One or more FPGAs used



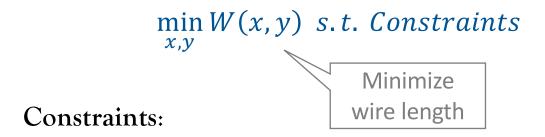
FPGA DESIGN IMPLEMENTATION



FPGA PLACEMENT

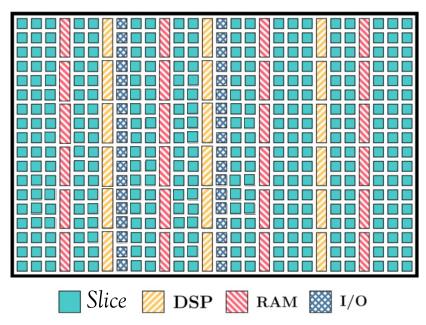


Placement Problem Formulation



- Density: No overlaps between instances,
- Routability: Design is routable,
- Congestion: Ensure placement is not congested etc.,

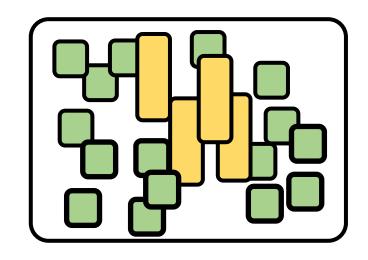
A simplified Xilinx Ultrascale architecture



Source: Xilinx; ISPD'2016 Contest

GLOBAL PLACEMENT

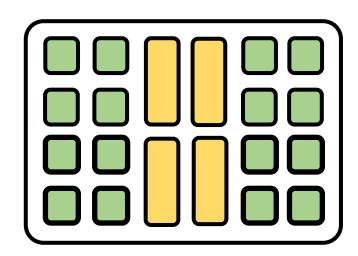
- Obtain rough legal locations for instances
- Optimization problem can be formulated as:
 - Simulated Annealing
 - Quadratic Placement
 - Non-linear Placement

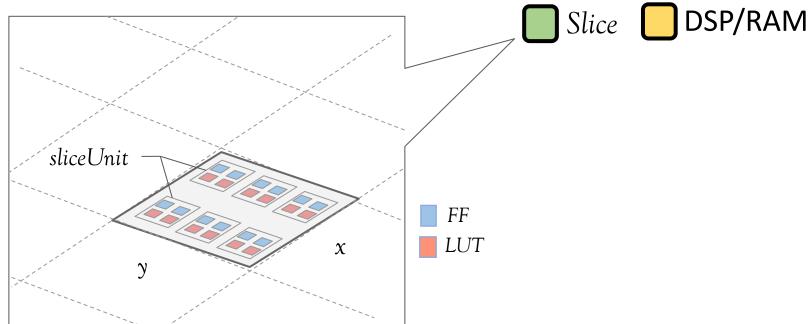


LEGALIZATION

Assign instances to respective sites on the FPGA

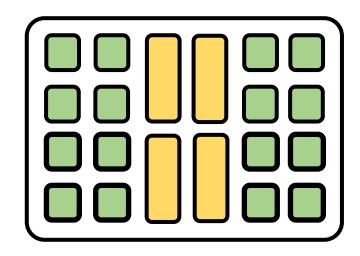
Look-up tables (LUTs) & FlipFlops (FFs) are packed



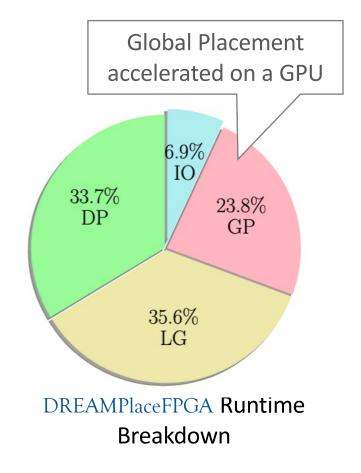


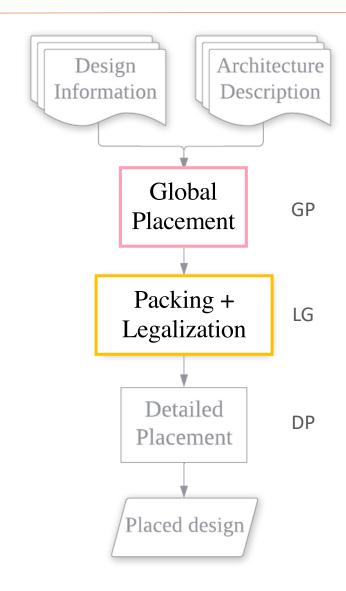
DETAILED PLACEMENT

- Further refine the legalized placement
- Formulated as
 - Independent Set Matching
 - Dynamic Programming
 - Hill Climbing, etc.



FPGA PLACEMENT ACCELERATION





EXISTING PACKER-LEGALIZER ACCELERATION SCHEMES

- Multi-threaded CPU [w.Li+, ICCAD'2019; w.Li+, TCAD'2017; L.Singhal+, DAC'2017]
 - Scalable consensus-based algorithms

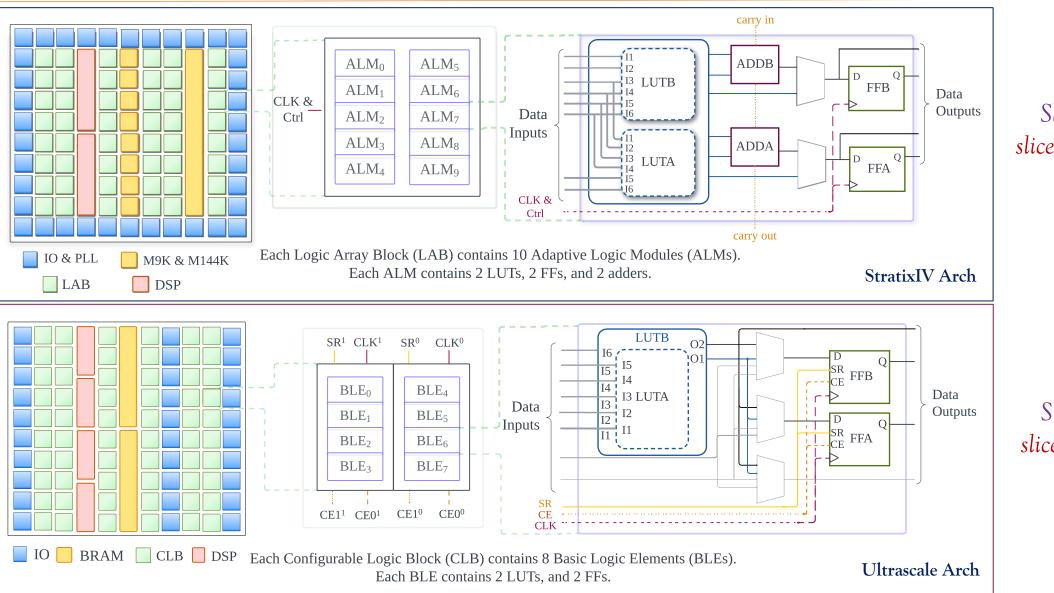
Other acceleration schemes not explored for FPGA packing-legalization!

- Use of CPU-GPU systems in ASIC LG Acceleration [H.Yang+, DATE'2022]
 - 2x 4x speedup observed
 - In ASIC, Legalization accounts for ~5% of total placement runtime [Y.Lin+, TCAD'2020]

EXISTING PACKER-LEGALIZER ACCELERATION SCHEMES

- Multi-threaded CPU [w.Li+, ICCAD'2019; w.Li+, TCAD'2017; L.Singhal+, DAC'2017]
 - Scalable consensus-based algorithms
 - Oth In FPGAs, LG takes ~35% of total placement runtime
- Accelerating on GPUs can help reduce runtime!
 - Zx 4x speedup observed
 - In ASIC, Legalization accounts for ~5% of total placement runtime [Y.Lin+, TCAD'2020]

PACKING RULES IN FPGA ARCHITECTURES



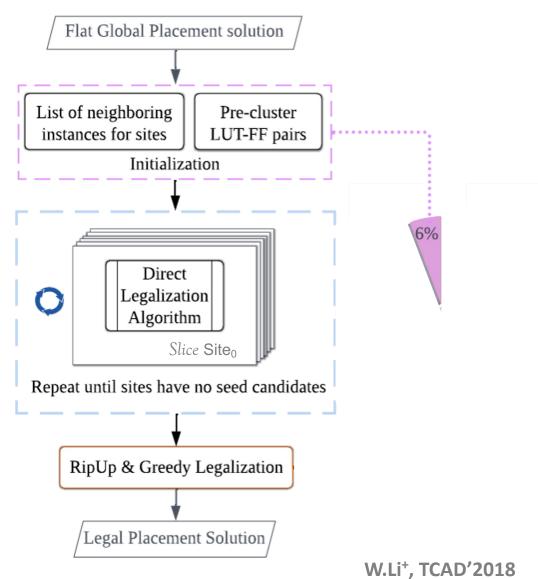
 $Slice \longleftrightarrow LAB$ $slice Unit \longleftrightarrow ALM$

 $Slice \longleftrightarrow CLB$ $sliceUnit \longleftrightarrow BLE$

SIMULTANEOUS PACK-LEGALIZE ALGORITHM

- Site centric algorithm ⇒ Scalable
- LUTs and FFs in a flat global placement are packed and legalized

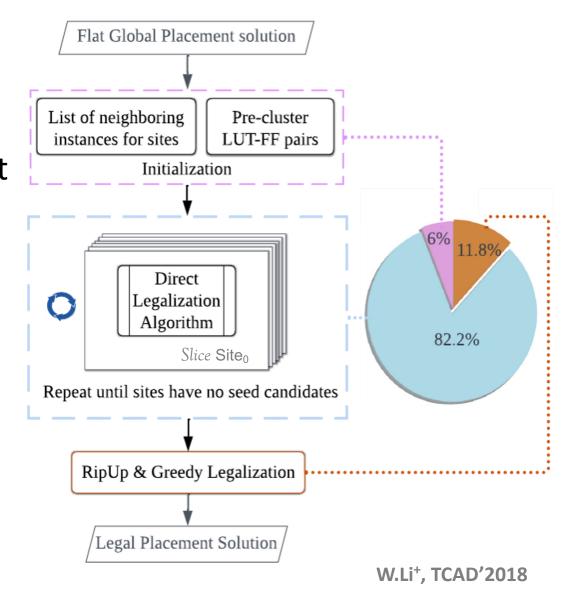
- Initialization:
 - Precluster LUT with FFs in fanout
 - Get neighbor instances of each site



SIMULTANEOUS PACK-LEGALIZE ALGORITHM

- Direct Legalization (DL) Algorithm:
 - Best Candidate for GPU Acceleration
 - A Candidate: cluster of LUTs + FFs that
 can be placed in a Slice
 - Each site runs in parallel
 - > 90% LUTs + FFs are legalized
- RipUp & Greedy Legalization:

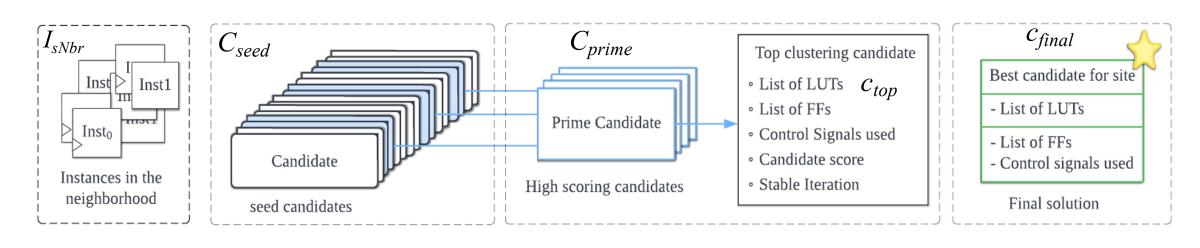
Handles the remaining instances



DREAMPLACEFPGA-PL

- First attempt to accelerate the pack-legalize stage on GPUs
- Accelerate the Direct Legalization (DL) algorithm
 - Revise runtime-critical portions in the algorithm
 - Employ task scheduling for even load across GPU threads
- Contribute to FPGA open-source ecosystem

FLAT DATASTRUCTURES FOR EACH SLICE



 $I_{
m sNbr}$: Set of all the neighboring instances of $Slice\ s$ within a distance D

Nbr: Instances currently considered by the Slice - a subset of I_{sNbr}

 C_{seed} : Set of cluster candidates for Slice s

 C_{prime} : Top 10 candidates in C_{seed} based on score (HPWL + internal nets)

 $c_{\it top}$: Best candidate in $C_{\it prime}$

 c_{final} : Final solution for Slice s

THE DIRECT LEGALIZATION (DL) ALGORITHM

Run in parallel for all Slice sites

Algorithm 1 The Direct Legalization Algorithm

```
1: while |I_{sNbr}| > 0 or |C_{seed}| > 0 do
       if c_{top} is stable for T iterations then
           c_{final} = c_{top}
           Clear all seed candidates in C_{seed}
           Add c_{top} to C_{seed} as seed candidate
 5:
       else
           Remove invalid candidates in C_{seed} and C_{prime}
7:
       end if
       Remove committed neighbor instances in Nbr and I_{sNbr}
       if |Nbr| < MinNbrCount then
10:
           add next group of neighbors from I_{sNbr} to Nbr;
11:
       end if
12:
       Create new seed candidates from C_{seed} and Nbr
13:
       Update best site for instances in c_{top}
14:
15: end while
```

THE DIRECT LEGALIZATION (DL) ALGORITHM

Run in parallel for all Slice sites

```
Algorithm 1 The Direct Legalization Algorithm
1: while |I_{sNbr}| > 0 or |C_{seed}| > 0 do
       if c_{top} is stable for T iterations then
           c_{final} = c_{top}
           Clear all seed candidates in C_{seed}
           Add c_{top} to C_{seed} as seed candidate
 5:
       else
           Remove invalid candidates in C_{seed} and C_{prime}
       end if
       Remove committed neighbor instances in Nbr and I_{sNbr}
                                                                      A naive GPU Implementation
       if |Nbr| < MinNbrCount then
10:
                                                                       is > 3x slower than 8T CPU
           add next group of neighbors from I_{sNbr} to Nbr;
11:
       end if
12:
                                                                   GPU Runtime
       Create new seed candidates from C_{seed} and Nbr
13:
                                                                 Bottleneck (RB)!
       Update best site for instances in c_{top}
15: end while
```

RB1: New Cluster Candidate Creation

- ullet New candidates = Neighbor Instances in Nbr imes candidates in C_{seed}
- \bullet Slice sites can have different | Nbr | and | C_{seed} | \Rightarrow Exploration space differs!
 - Some Slice sites can have no new candidates
 - Some could have 100s-1000s of new candidates => Load imbalance!

RB1: MITIGATION SCHEMES

(1) Staggered addition of new neighbor instances

- Default: Group of neighbors added
- Stagger $|\varsigma|$: Max of $|\varsigma|$ neighbors added

 $|\varsigma|$: Max number of LUTs + FFs in Slice

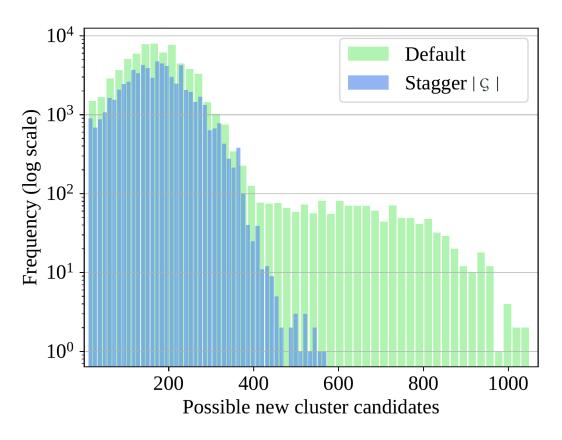
For example, consider a Slice with 100 neighbor instances

All neighbor instances are at a distance 3 from Slice.

Let $|\varsigma| = 32$. New neighbor addition is as follows:

Default: {0, 0, 100}

Stagger | \(\mathbf{\sigma} \) |: \(\{32, 32, 32, 4\} \)



Histogram of new candidate search space of all the sites at the fourth DL algorithm iteration in FPGA12.

RB1: MITIGATION SCHEMES

- (1) Staggered addition of new neighbor instances
- **(2) Scheduling** of *Slice* sites
 - ullet Order Slice sites in decreasing order of new candidate exploration space
 - \bullet Work on Slice sites with similar load at any given time

RB1: MITIGATION SCHEMES

- (1) Staggered addition of new neighbor instances
- **(2) Scheduling** of *Slice* sites
- (3) Restrict total new candidates explored
 - \bullet Limit the total | Nbr | x | C_{seed} | combinations explored
 - Could degrade placement quality as potential candidates could be discarded

RB2: Instance Best Site Assignment

- ullet Instances in c_{top} select the best Slice site based on:
 - HPWL improvement for instance to move to Slice
 - Possiblity of internal nets absorbed in a Slice
- \bullet Tie-break if multiple *Slice* sites have same score low site identifier
- ◆ Each Instance selects best site sequentially ⇒ Stalls GPU threads!
- Changes to instance datastructure increases memory footprint significantly

RB2: MITIGATION

- Minimize sequential portion and include post-processing step
 - Slower on CPU but much faster on GPU
 - Solution is different but there is no loss of placement quality!

```
set_lock(); //Begin sequential portion
if (site.score == inst.bestScore)
    //Tie-break condition
    if (site.id < inst.bestSiteId)</pre>
        inst.bestSiteId = site.id;
else if (site.score > inst.bestScore)
    inst.bestSiteId = site.id;
    inst.bestScore = site.score;
unset_lock(); //End sequential portion
```

```
//Sequential assignment of
// instance best score
atomicMax(&inst.bestScore, site.score);
//Parallel post-processing
//Get instance bestSiteId
// based on bestScore
for (     sites in neighborhood of inst)
    if (site.score == inst.bestScore &&
        site.top.candidate contains inst &&
        site.id < inst.bestSiteId) //tie-break</pre>
        inst.bestSiteId = site.id;
//End of parallel post-processing
```

EXPERIMENTAL SETUP

• **Design suite:** ISPD'2016 benchmarks

Placers: elfPlace

• Router: Xilinx Vivado v2015.4

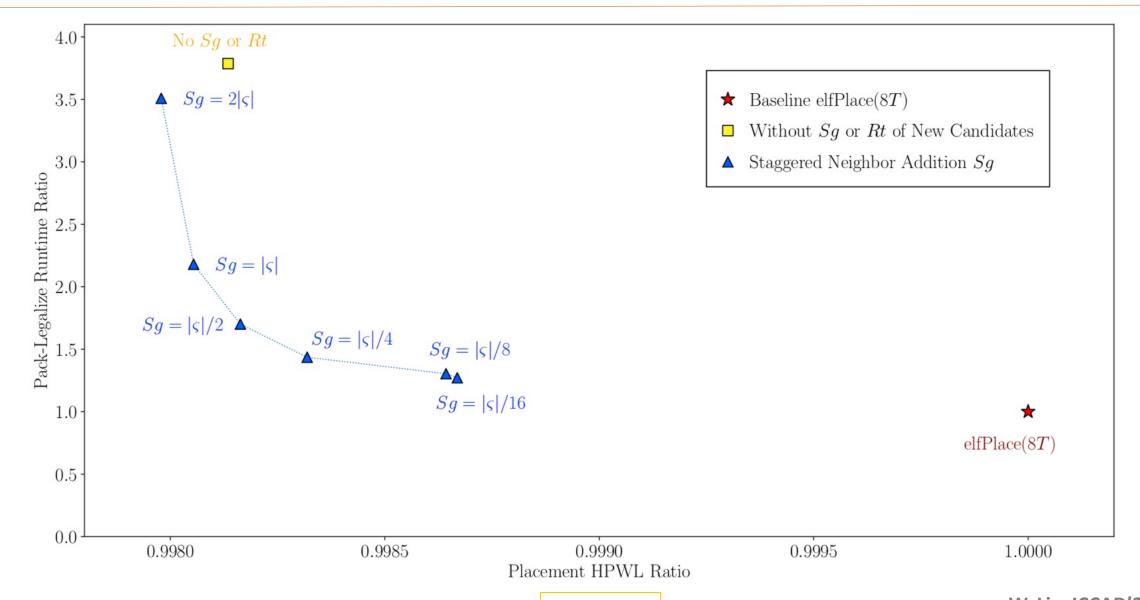
- Comparison Metrics
 - Placement HPWL, Routed Wirelength
 - Placement Runtime

 Design	#IO	#RAM	#DSP	#LUT	#FF
————	π1Ο	πICTIVI	π1001	#LU 1	π11
FPGA01	156	0	0	50k	55k
FPGA02	156	100	100	100k	66k
FPGA03	406	600	500	250k	170k
FPGA04	406	600	500	250k	172k
FPGA05	406	600	500	250k	174k
FPGA06	606	1000	600	350k	352k
FPGA07	606	1000	600	350k	355k
FPGA08	406	600	500	500k	216k
FPGA09	606	1000	600	500k	366k
FPGA10	606	1000	600	350k	600k
FPGA11	606	1000	400	480k	363k
FPGA12	406	600	500	500k	602k
Resources	2580	1728	768	538k	1075k

^{*}DP run on elfPlace

STAGGERED NEIGHBOR ADDITION

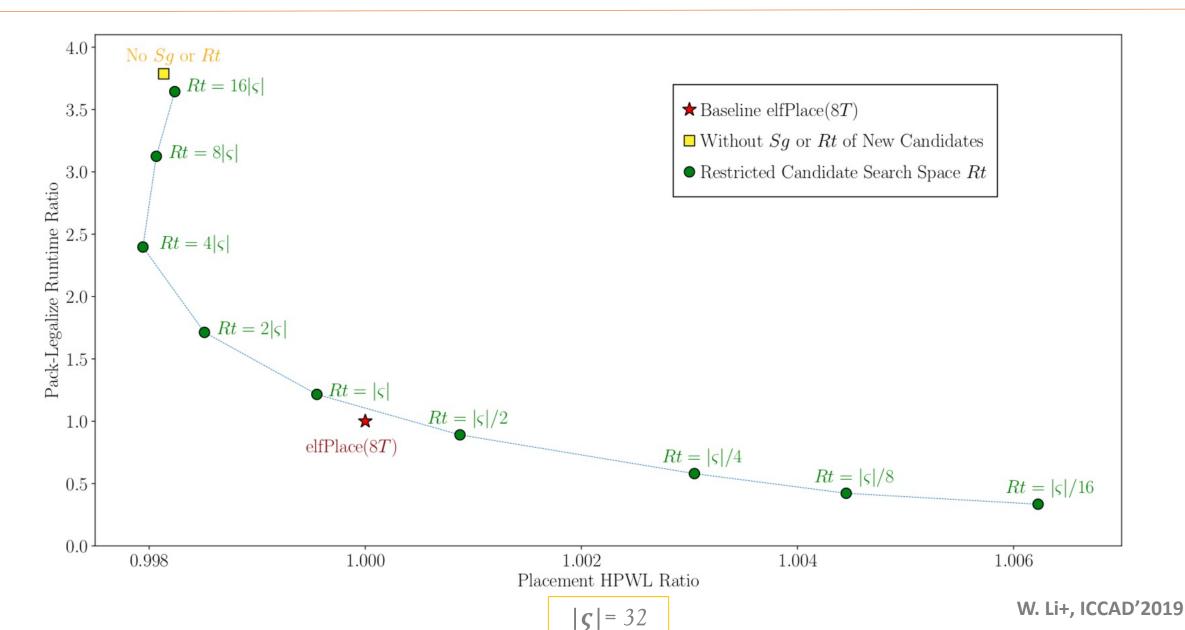
GPU: NVIDIA TITAN Xp (Pascal)



GDI I

GPU: NVIDIA TITAN Xp (Pascal)

RESTRICTED NEW CANDIDATE EXPLORATION



SG + RT VARIANTS

Two variants that produce the best runtime vs QoR tradeoff considered are:

$$\bullet$$
 V1: Sg = $|\varsigma|/16 + Rt = |\varsigma|/8$

- Aggressive limits on Staggered neighbor addition and Restricted search space
- Achieves the best runtime speedup with a slight degradation of QoR

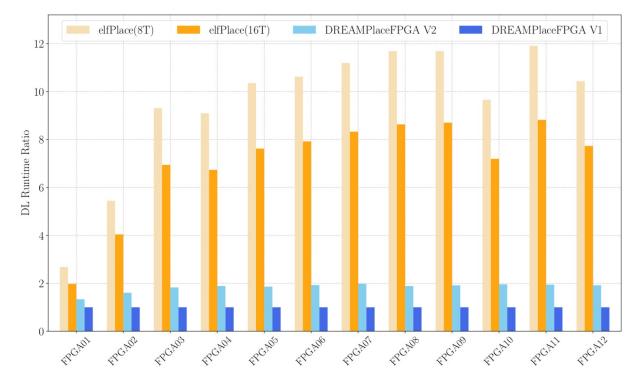
•
$$V2: Sg = |\varsigma|/8 + Rt = |\varsigma|/4$$

- Less aggressive than V1
- Has lower QoR degradation and runtime improvement

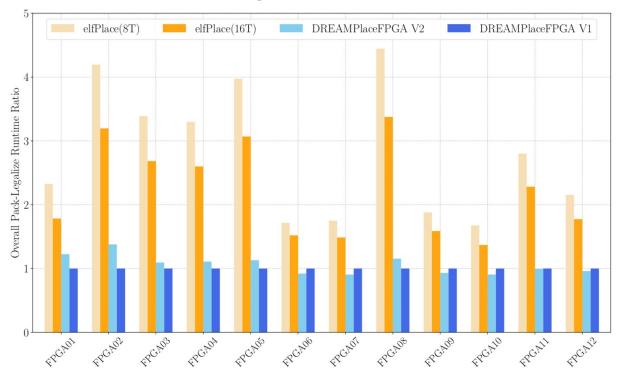
RUNTIME IMPROVEMENTS

GPU: NVIDIA TITAN Xp (Pascal)

Direct Legalization (DL) Runtime



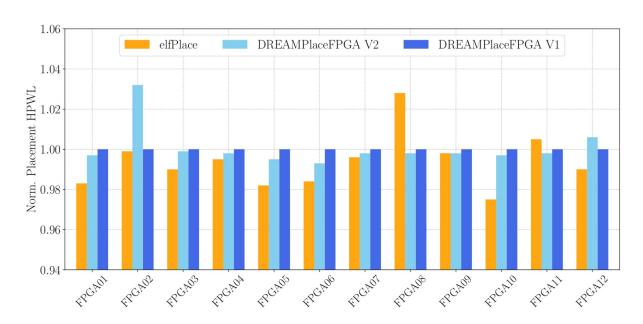
Pack-Legalize (LG) Runtime



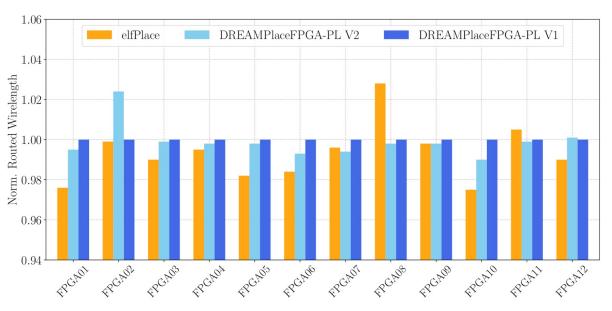
Metric	elfPlace(8T)	elfPlace(16T)	Ours V2	Ours V1
Init	0.97	0.97	0.92	1.00
Direct Legalization (DL)	9.51	7.05	1.84	1.00
RipUp & Greedy LG (RG)	0.79	0.76	0.74	1.00
Total LG	2.80	2.23	1.06	1.00

IMPACT ON PLACEMENT QUALITY





Routed Wirelength



Metric	elfPlace	Ours V2	Ours V1
Placement HPWL	0.994	0.999	1.000
Routed Wirelength	0.991	0.997	1.000

W. Li+, ICCAD'2019

DREAMPLACEFPGA-PL SUMMARY

- First attempt to accelerate FPGA packer-legalizer on GPU
- Revise the direct legalization algorithm => 7x speedup than 16T CPU
- DREAMPlaceFPGA-PL
 - 2.23x faster LG than 16T CPU
 - +0.6% placement HPWL and +0.9% routed WL
- With GP and LG accelerated on GPU, placement is 2.38x faster than 16T CPU

THANK YOU