

GSOHC:Global Synchronization Optimization for Heterogeneous Computing



भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad

Soumik Kumar Basu

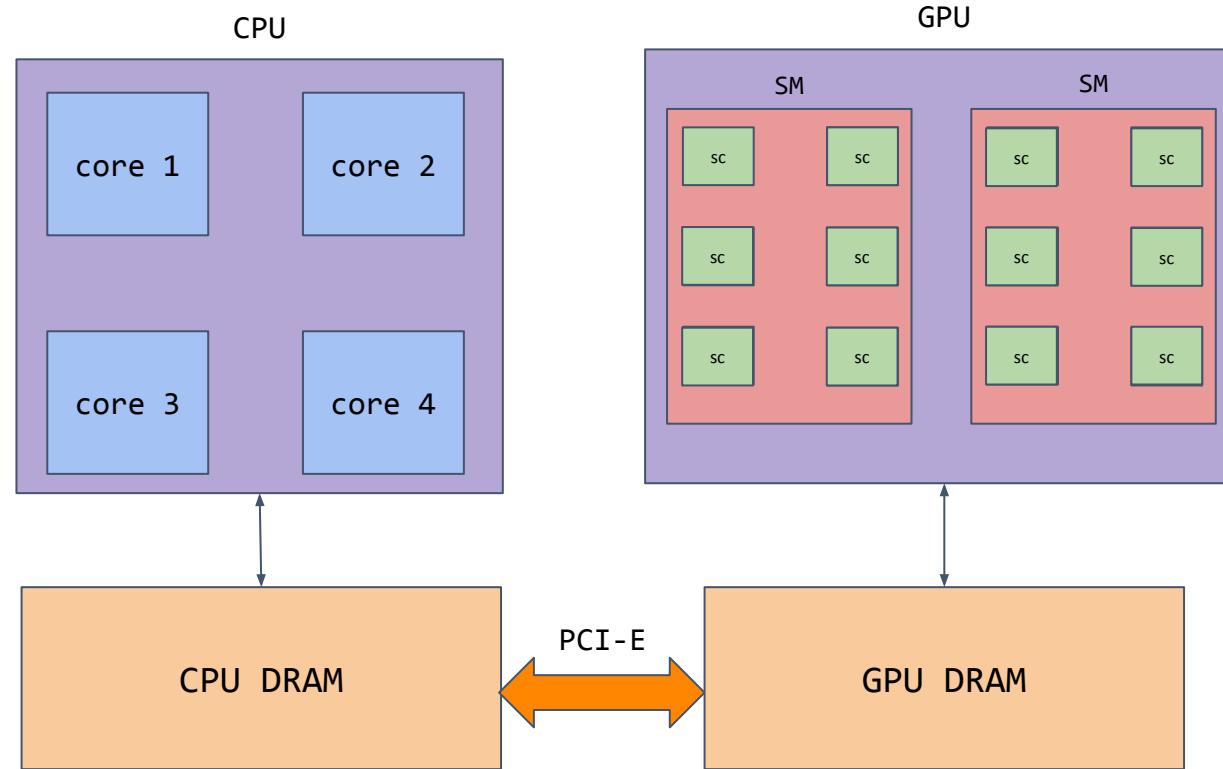
Jyothi Vedurada

Department of Computer Science and Engineering

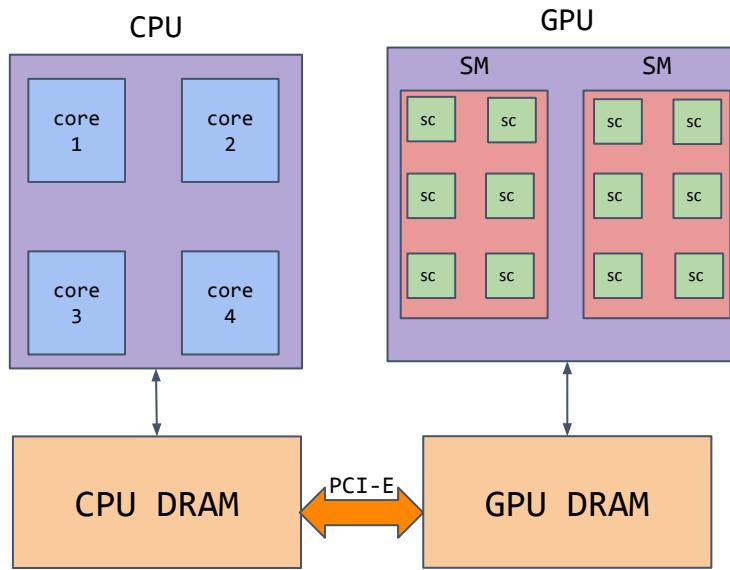
IIT Hyderabad



Execution of a heterogeneous program



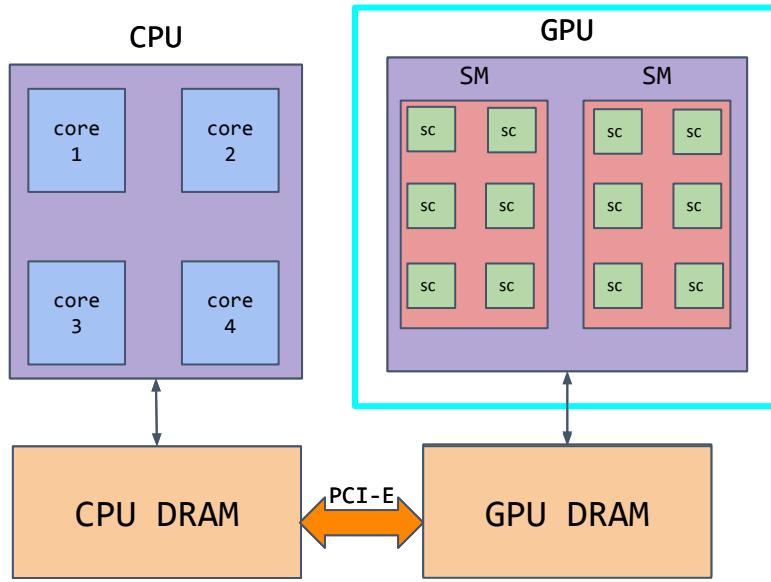
A typical execution of CUDA programs



```
16. void foo(int *x, int size2)
17. {
18. ...
19. if(k < y_points)
20.     x[k] = x[0] + y_points;
21. ...
22. }
```

```
1. void main(){
2. ...
3. kernel2<<<grid2, block2>>>(d_u,...);/*Device-side computations*/
4.
5. cudaMemcpy(v, d_u, size2, cudaMemcpyDeviceToHost);
6.
7. cudaDeviceSynchronize();/*Sync*/
/*Host-side computations*/
8. ...
9.
10. foo(v, size2);
11. ...
12. for(int i=0; i<size2; i++)
13.     v[i] = v[i] + x_points;
14. ...
15. }
```

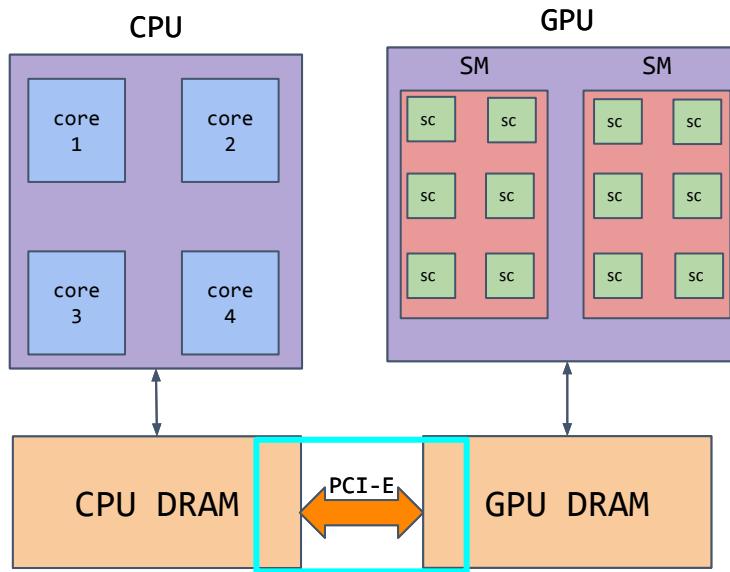
A typical execution of CUDA programs



```
16. void foo(int *x, int size2)
17. {
18. ...
19. if(k < y_points)
20.     x[k] = x[0] + y_points;
21. ...
22. }
```

```
1. void main(){
2. ...
3. kernel2<<<grid2, block2>>>(d_u,...);/*Device-side computations*/
4.
5. cudaMemcpy(v, d_u, size2, cudaMemcpyDeviceToHost);
6.
7. cudaDeviceSynchronize();/*Sync*/
8. /*Host-side computations*/
9.
10. ...
11. foo(v, size2);
12. ...
13. for(int i=0; i<size2; i++)
14.     v[i] = v[i] + x_points;
15. }
```

A typical execution of CUDA programs



```
16. void foo(int *x, int size2)
17. {
18. ...
19. if(k < y_points)
20.     x[k] = x[0] + y_points;
21. ...
22. }
```

```
1. void main(){
2. ...
3. kernel2<<<grid2, block2>>>(d_u,...);/*Device-side computations*/
4.
5. cudaMemcpy(v, d_u, size2, cudaMemcpyDeviceToHost);
6.
7. cudaDeviceSynchronize();/*Sync*/
8. /*Host-side computations*/
9.
10. ...
11. foo(v, size2);
12. ...
13. for(int i=0; i<size2; i++)
14.     v[i] = v[i] + x_points;
15. }
```

A brief introduction of CUDA programs

```
1. void main(){
2. ...
3. kernel2<<<grid2, block2>>>(d_u,...);/*Device-side computations*/
4.
5. cudaMemcpy(v, d_u, size2, cudaMemcpyDeviceToHost);
6.
7. cudaDeviceSynchronize();/*Sync*/
8. /*Host-side computations*/
9. ...
10. foo(v, size2);
11. ...
12. for(int i=0; i<size2; i++)
13.     v[i] = v[i] + x_points;
14. ...
15. }
```

A brief introduction of CUDA programs

```
1. void main(){
2. ...
3. kernel2<<<grid2, block2>>>(d_u,...);/*Device-side computations*/
4.
5. cudaMemcpy(v, d_u, size2, cudaMemcpyDeviceToHost);
6.
7. cudaDeviceSynchronize();/*Sync*/
8. /*Host-side computations*/
9. ...
10. foo(v, size2);
11. ...
12. for(int i=0; i<size2; i++)
13.     v[i] = v[i] + x_points;
14. ...
15. }
```

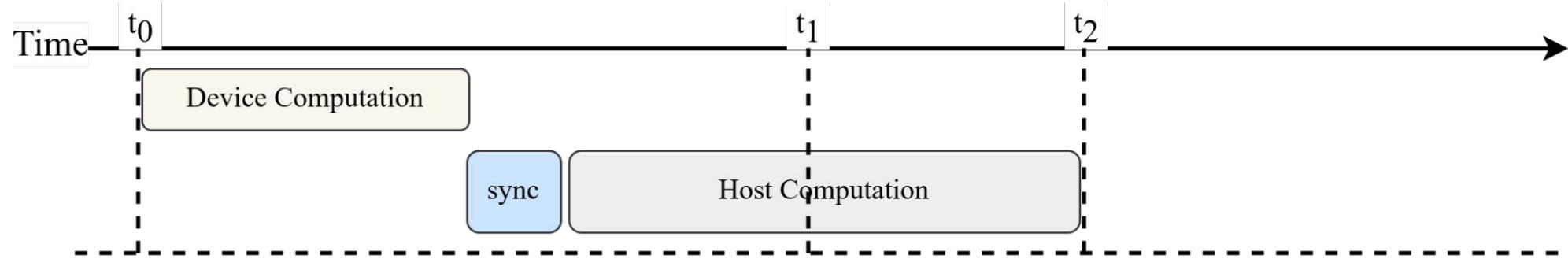
A brief introduction of CUDA programs

Blocking
Calls!!

```
16. void foo(int *x, int size2)
17. {
18.   ...
19.   if(k < y_points)
20.     x[k] = x[0] + y_points;
21.   ...
22. }
```

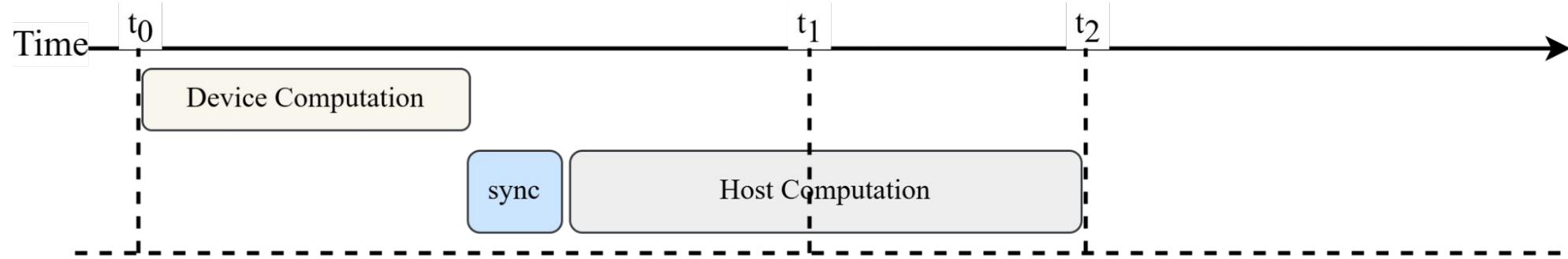
```
1. void main(){
2. ...
3. kernel2<<<grid2, block2>>>(d_u,...);/*Device-side computations*/
4.
5. cudaMemcpy(v, d_u, size2, cudaMemcpyDeviceToHost);
6.
7. cudaDeviceSynchronize();/*Sync*/
8. /*Host-side computations*/
9. ...
10. foo(v, size2);
11. ...
12. for(int i=0; i<size2; i++)
13.   v[i] = v[i] + x_points;
14. ...
15. }
```

Problem: Poor placement of synchronization statements



Total time taken = $t_2 - t_0$

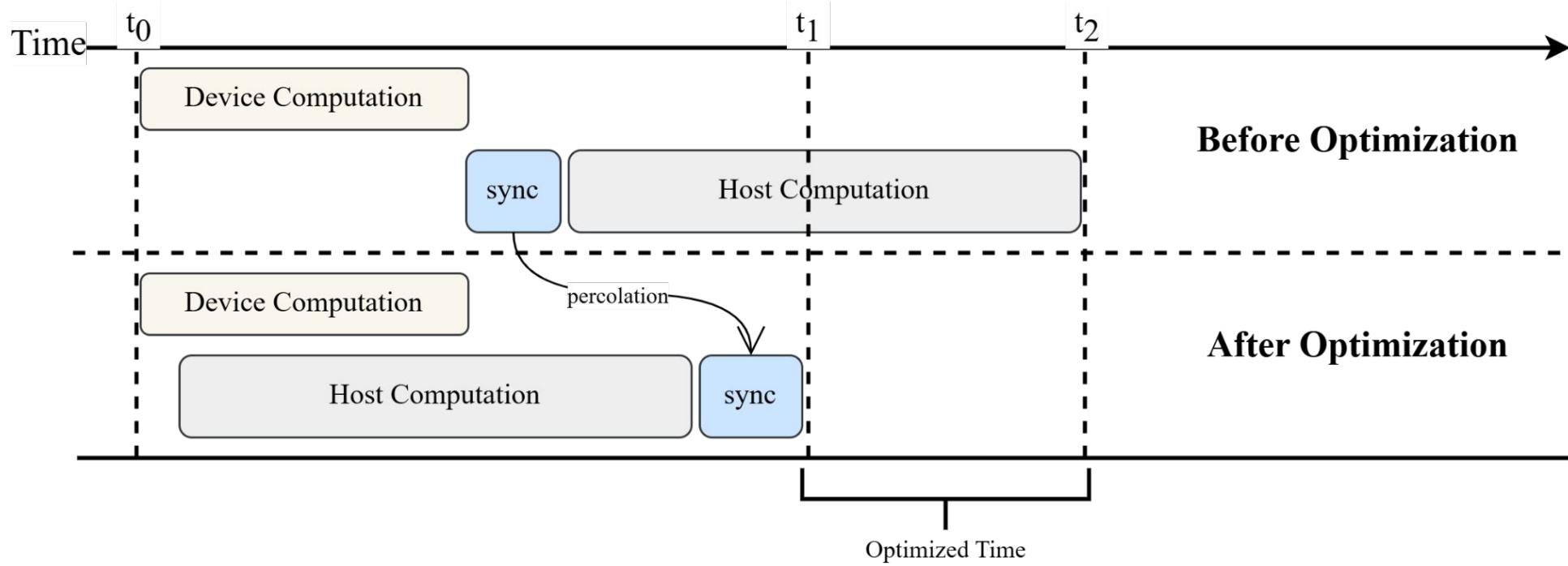
Problem: Poor placement of synchronization statements



$$\text{Total time taken} = t_2 - t_0$$

How can we improve
such a program ?

Solution: Hetero-sync motion



$$\text{Total time taken} = t_1 - t_0$$

Hetero-sync motion opportunity: Graph Coloring (Case Study)

```
1. void divideWork(...){  
2.     if(value < size) {  
3.         for(int i = 0; i < value; ++i) {...}  
4.     }  
5.     else {  
6.         int portion = value / size;  
7.         for(int i = 0; i < size; ++i) {  
8.             h_head[i] = i * portion;  
9.             h_tail[i] = (i + 1) * portion;  
10.        }  
11.        h_tail[size - 1] = value;  
12.    }  
13.    cudaMemcpy(d_head, h_head, ...,  
14.    cudaMemcpyHostToDevice);  
15.    cudaMemcpy(d_tail, h_tail, ...,  
16.    cudaMemcpyHostToDevice);  
17. }
```

```
1. void main(){...  
2.     while(!(*h_complete)) {  
3.         divideWork(h_head, h_tail, d_head, d_tail, NBLOCKS, v);  
4.         assignColorsKernel<<<...>>> (d_head, ...);  
5.     }  
6.     cudaDeviceSynchronize();  
7.     *h_complete = true;  
8.     cudaMemcpy(d_complete, ..., cudaMemcpyHostToDevice);  
9.     divideWork(...);  
10.    detectConflictsKernel<<<...>>> (... , d_complete, ...);  
11.    cudaDeviceSynchronize();  
12.    cudaMemcpy(..., d_complete, ..., cudaMemcpyDeviceToHost);  
13.    divideWork(...);  
14.    forbidColorsKernel<<<...>>> (d_head, ...);  
15.    cudaDeviceSynchronize();}  
16.    ...  
17. }
```

Hetero-sync motion opportunity: Graph Coloring (Case Study)

```
1. void divideWork(...){  
2.     if(value < size) {  
3.         for(int i = 0; i < value; ++i) {...}  
4.     }  
5.     else {  
6.         int portion = value / size;  
7.         for(int i = 0; i < size; ++i) {  
8.             h_head[i] = i * portion;  
9.             h_tail[i] = (i + 1) * portion;  
10.        }  
11.        h_tail[size - 1] = value;  
12.    }  
13.    cudaMemcpy(d_head, h_head, ...,  
14.    cudaMemcpyHostToDevice);  
15.    cudaMemcpy(d_tail, h_tail, ...,  
16.    cudaMemcpyHostToDevice);  
17. }
```

```
1. void main(){...  
2.     while(!(*h_complete)) {  
3.         divideWork(h_head, h_tail, d_head, d_tail, NBLOCKS, v);  
4.         assignColorsKernel<<<...>>> (d_head, ...);  
5.     }  
6.     cudaDeviceSynchronize();  
7.     *h_complete = true;  
8.     cudaMemcpy(d_complete, ..., cudaMemcpyHostToDevice);  
9.     divideWork(...);  
10.    detectConflictsKernel<<<...>>> (..., d_complete, ...);  
11.    cudaDeviceSynchronize();  
12.    cudaMemcpy(..., d_complete, ..., cudaMemcpyDeviceToHost);  
13.    divideWork(...);  
14.    forbidColorsKernel<<<...>>> (d_head, ...);  
15.    cudaDeviceSynchronize();}  
16.    ...  
17. }
```

Hetero-sync motion opportunity: Graph Coloring (Case Study)

```
1. void divideWork(...){  
2.     if(value < size) {  
3.         for(int i = 0; i < value; ++i) {...}  
4.     }  
5.     else {  
6.         int portion = value / size;  
7.         for(int i = 0; i < size; ++i) {  
8.             h_head[i] = i * portion;  
9.             h_tail[i] = (i + 1) * portion;  
10.        }  
11.        h_tail[size - 1] = value;  
12.    }  
13.    cudaMemcpy(d_head, h_head, ...,  
14.    cudaMemcpyHostToDevice);  
15.    cudaMemcpy(d_tail, h_tail, ...,  
16.    cudaMemcpyHostToDevice);  
17. }
```

```
1. void main(){...  
2.     while(!(*h_complete)) {  
3.         divideWork(h_head, h_tail, d_head, d_tail, NBLOCKS, v);  
4.         assignColorsKernel<<<...>>> (d_head, ...);  
5.         cudaDeviceSynchronize();  
6.         *h_complete = true;  
7.         cudaMemcpy(d_complete, ..., cudaMemcpyHostToDevice);  
8.         divideWork(...);  
9.         detectConflictsKernel<<<...>>> (... , d_complete, ...);  
10.        cudaDeviceSynchronize();  
11.        cudaMemcpy(..., d_complete, ..., cudaMemcpyDeviceToHost);  
12.        divideWork(...);  
13.        forbidColorsKernel<<<...>>> (d_head, ...);  
14.        cudaDeviceSynchronize();}  
15. }
```

Challenges of manually applying hetero-sync motion

Where to
Percolate ?

```
16. void foo(int *x, int size2)
17. {
18. ...
19. if(k < y_points)
20.     x[k] = x[0] + y_points;
21. ...
22. }
```

```
1. void main(){
2. ...
3. kernel2<<<grid2, block2>>>(d_u,...);/*Device-side computations*/
4.
5. cudaMemcpy(v, d_u, size2, cudaMemcpyDeviceToHost);
6.
7. cudaDeviceSynchronize();/*Sync*/
/*Host-side computations*/
8.
9. ...
10. foo(v, size2);
11. ...
12. for(int i=0; i<size2; i++)
13.     v[i] = v[i] + x_points;
14. ...
15. }
```

Challenges

```
16. void foo(int *x, int size2)
17. {
18.   ...
19.   if(k < y_points)
20.     x[k] = x[0] + y_points;
21.   ...
22. }
```

Challenge 1: Need to check for the data dependence while percolation to maintain program correctness

```
1. void main(){
2. ...
3. kernel2<<<grid2, block2>>>(d_u,...);/*Device-side computations*/
4.
5. cudaMemcpy(v, d_u, size2, cudaMemcpyDeviceToHost);
6.
7. cudaDeviceSynchronize();/*Sync*/
/*Host-side computations*/
8.
9. ...
10. foo(v, size2);
11. ...
12. for(int i=0; i<size2; i++)
13.   v[i] = v[i] + x_points;
14. ...
15. }
```

Challenges

```
16. void foo(int *x, int size2)
17. {
18.   ...
19.   if(k < y_points)
20.     x[k] = x[0] + y_points;
21.   ...
22. }
```

Challenge 1: Need to check for the data dependence while percolation to maintain program correctness

```
1. void main(){
2. ...
3. kernel2<<<grid2, block2>>>(d_u,...);/*Device-side computations*/
4.
5. cudaMemcpy(v, d_u, size2, cudaMemcpyDeviceToHost);
6.
7. cudaDeviceSynchronize();/*Sync*/
/*Host-side computations*/
8.
9. ...
10. foo(v, size2);
11. ...
12. for(int i=0; i<size2; i++)
13.   v[i] = v[i] + x_points;
14. ...
15. }
```

Challenges

```
16. void foo(int *x, int size2)
17. {
18.   ...
19.   if(k < y_points)
20.     x[k] = x[0] + y_points;
21.   ...
22. }
```

Challenge 2: Simply checking for data dependency may result in incorrectly percolation

```
1. void main(){
2. ...
3. kernel2<<<grid2, block2>>>(d_u,...);/*Device-side computations*/
4.
5. cudaMemcpy(v, d_u, size2, cudaMemcpyDeviceToHost);
6.
7. cudaDeviceSynchronize();/*Sync*/
/*Host-side computations*/
8.
9. ...
10. foo(v, size2);
11. ...
12. for(int i=0; i<size2; i++)
13.   v[i] = v[i] + x_points;
14. ...
15. }
```

Challenges

```
16. void foo(int *x, int size2)
17. {
18.   ...
19.   if(k < y_points)
20.     x[k] = x[0] + y_points;
21.   ...
22. }
```

Challenge 3: While percolation,
need to consider the farthest
target location

```
1. void main(){
2. ...
3. kernel2<<<grid2, block2>>>(d_u,...);/*Device-side computations*/
4.
5. cudaMemcpy(v, d_u, size2, cudaMemcpyDeviceToHost);
6.
7. cudaDeviceSynchronize();/*Sync*/
/*Host-side computations*/
8.
9. ...
10. foo(v, size2);
11. ...
12. for(int i=0; i<size2; i++)
13.   v[i] = v[i] + x_points;
14. ...
15. }
```

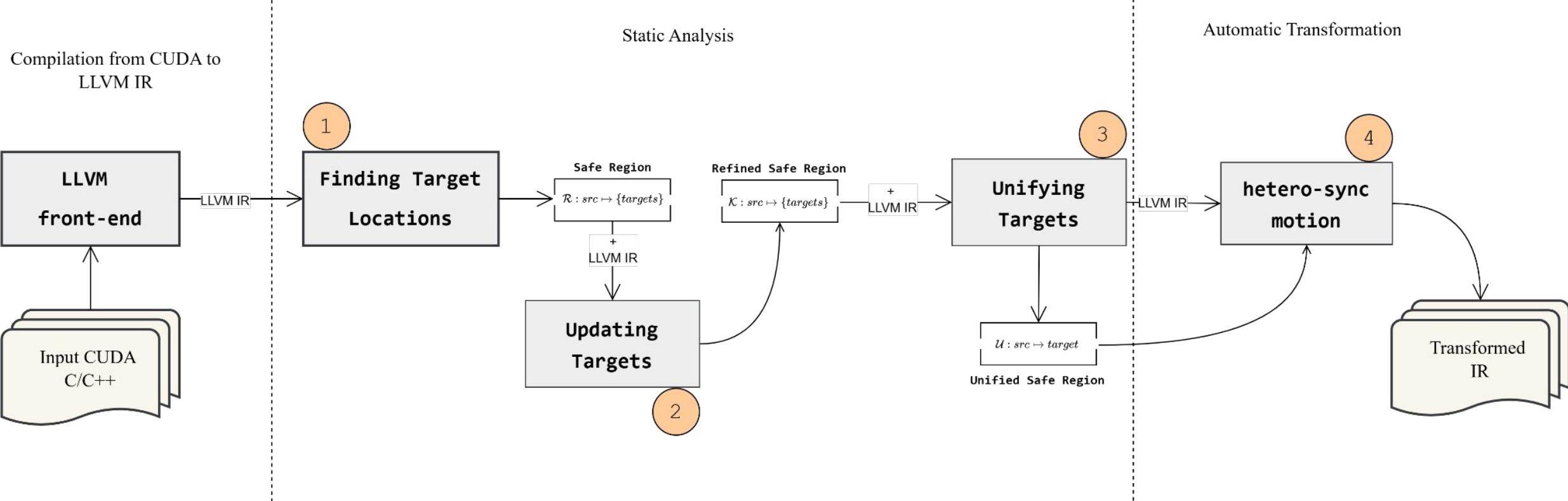
Challenges

```
16. void foo(int *x, int size2)
17. {
18.   ...
19.   if(k < y_points)
20.     x[k] = x[0] + y_points;
21.   ...
22. }
```

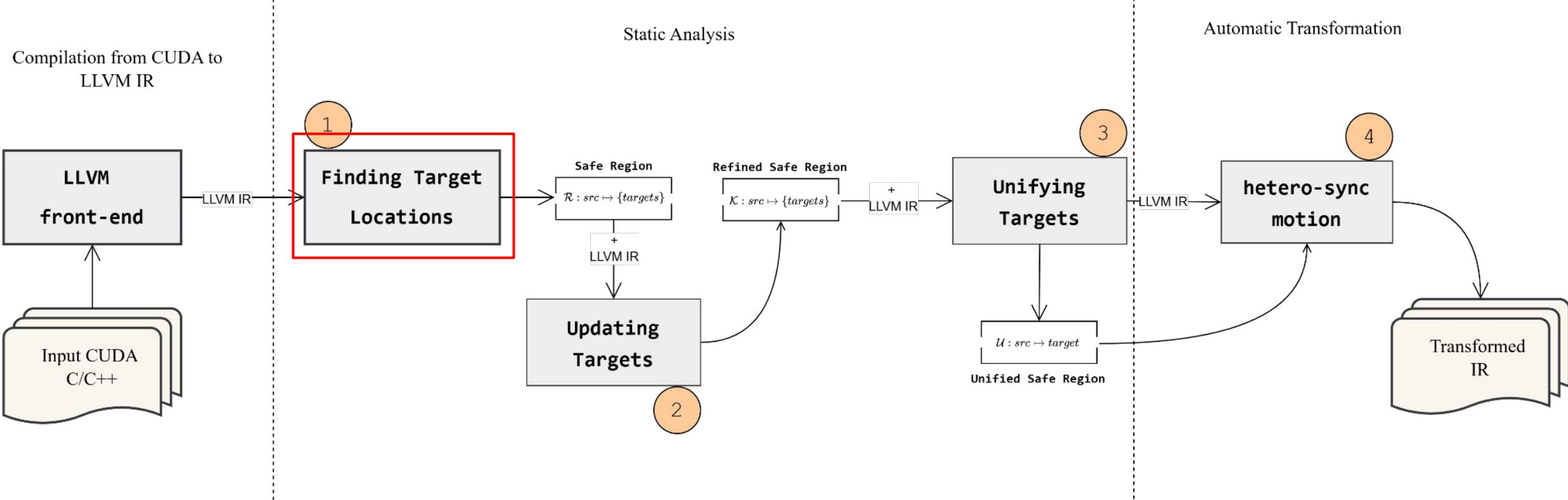
Challenge 4: A program may have multiple memcpy statements

```
1. void main(){
2. ...
3. kernel2<<<grid2, block2>>>(d_u,...);/*Device-side computations*/
4.
5. cudaMemcpy(v, d_u, size2, cudaMemcpyDeviceToHost);
6.
7. cudaDeviceSynchronize();/*Sync*/
/*Host-side computations*/
8.
9. ...
10. foo(v, size2);
11. ...
12. for(int i=0; i<size2; i++)
13.   v[i] = v[i] + x_points;
14. ...
15. }
```

Solution: GSOHC

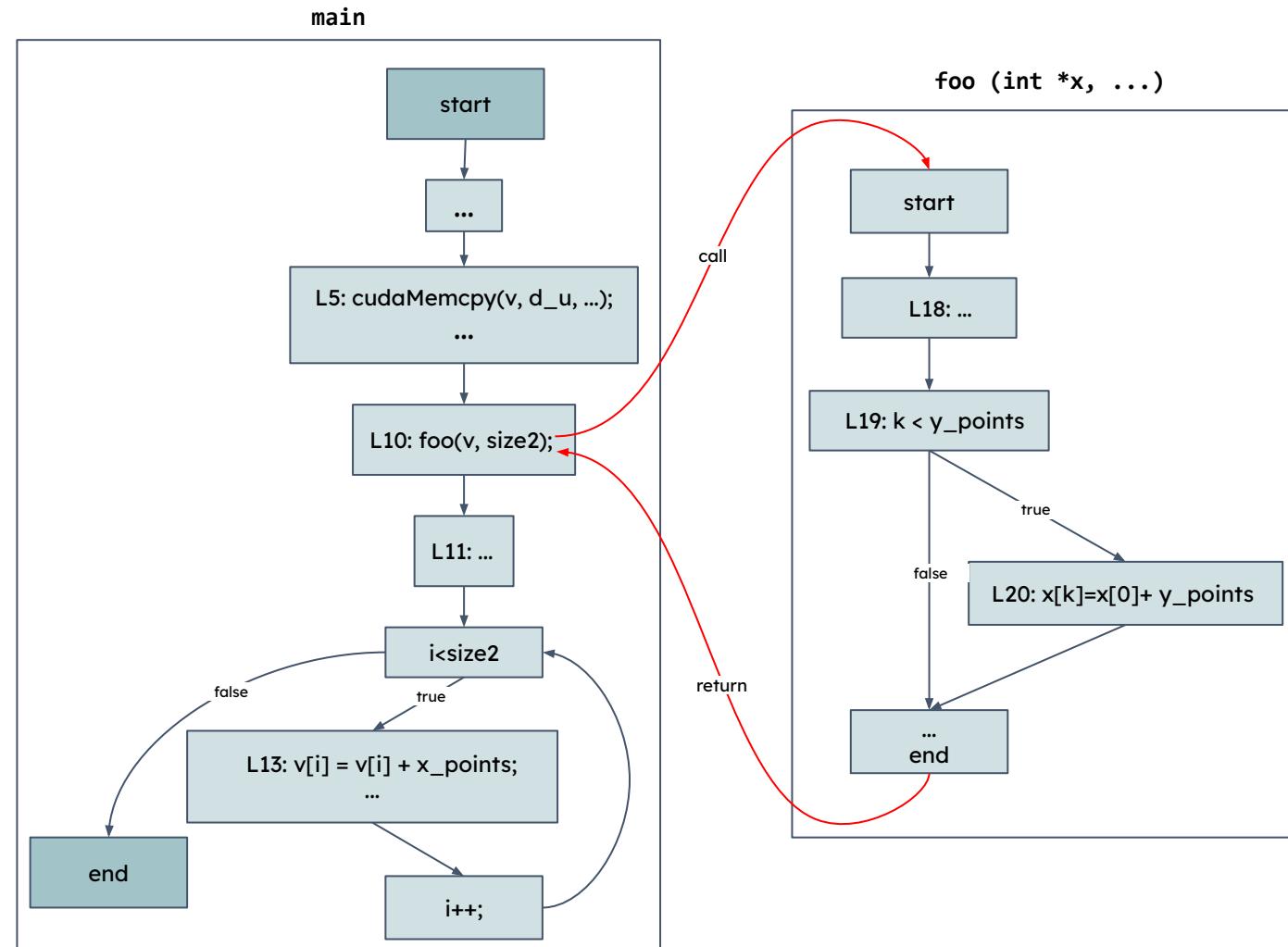


Step 1: Finding Target Locations

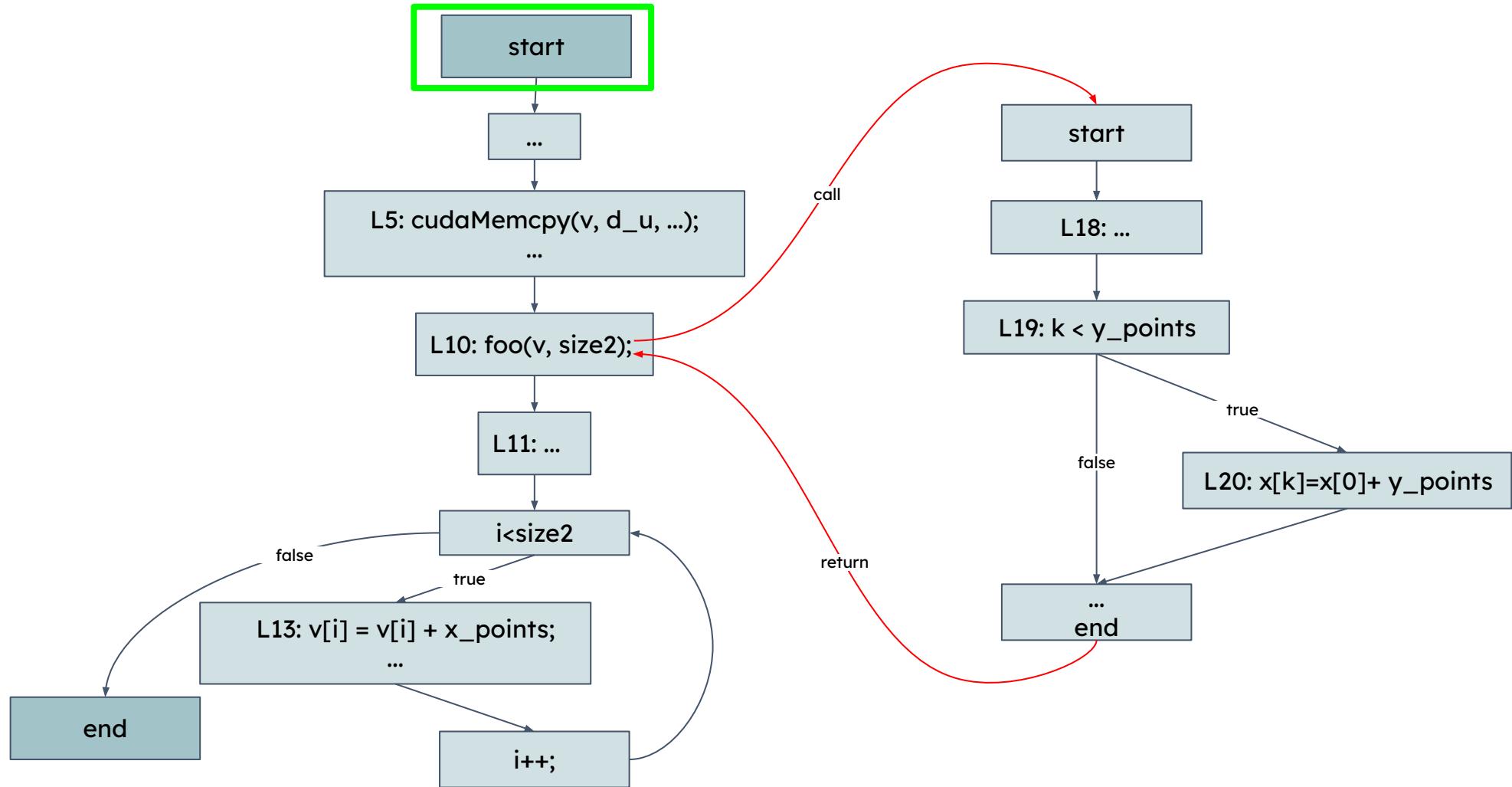


Step 1: Finding Target Locations

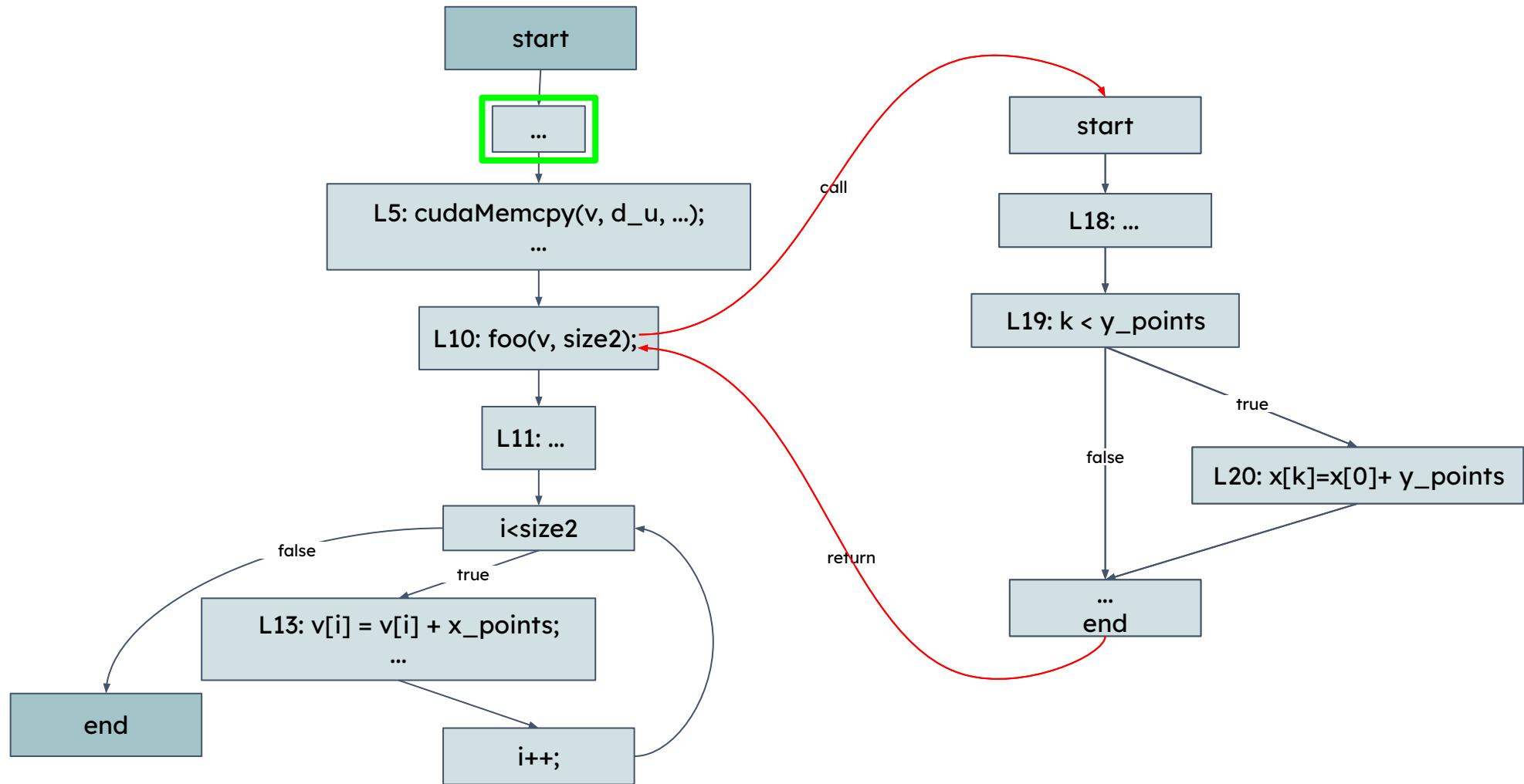
- Inter-procedural flow-sensitive, context-sensitive data-flow analysis
- Identifies target locations based on the data-dependence
- Handles multiple blocking calls at once



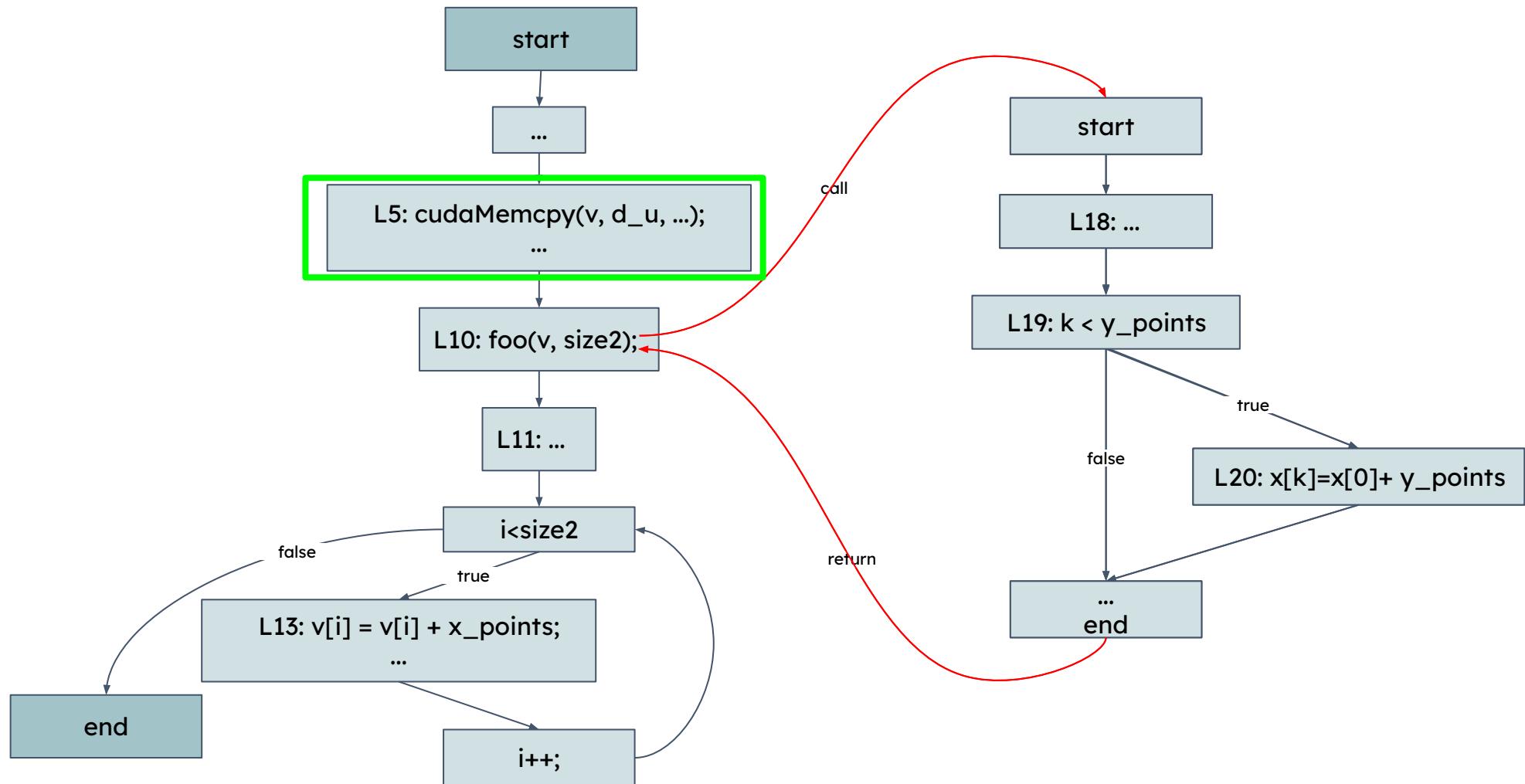
Step 1: Finding Target Locations (Contd.)



Step 1: Finding Target Locations (Contd.)

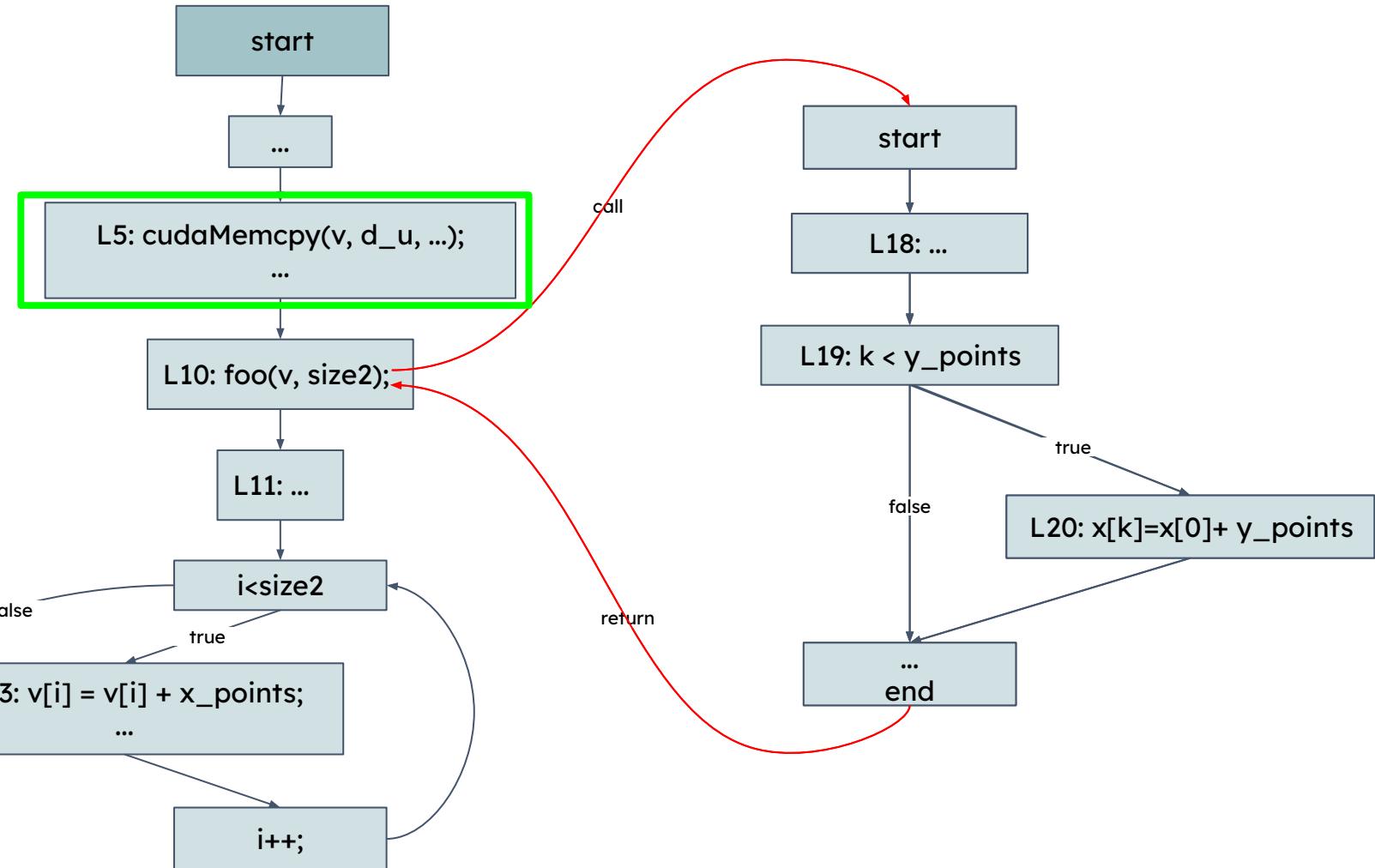


Step 1: Finding Target Locations (Contd.)

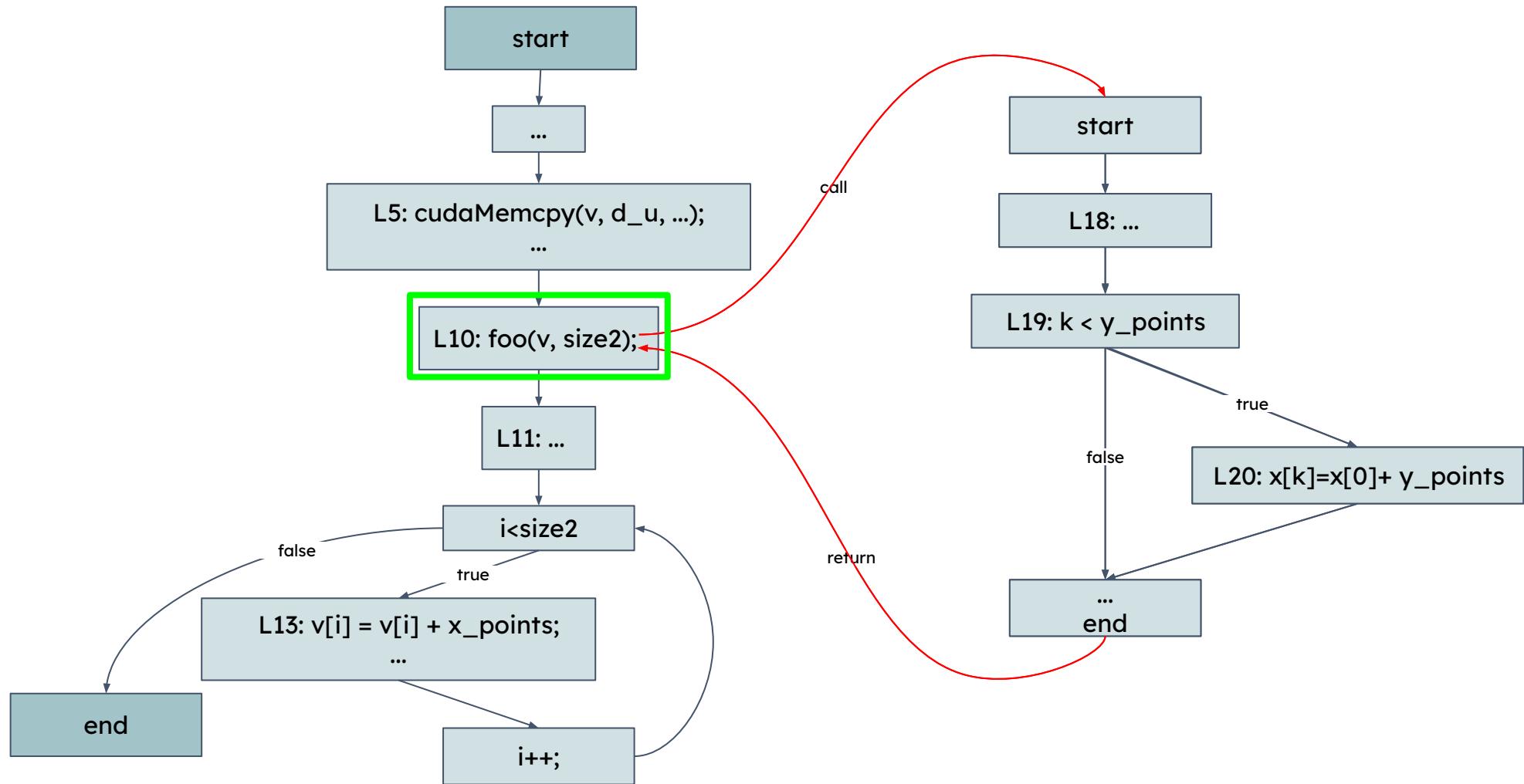


Step 1: Finding Target Locations (Contd.)

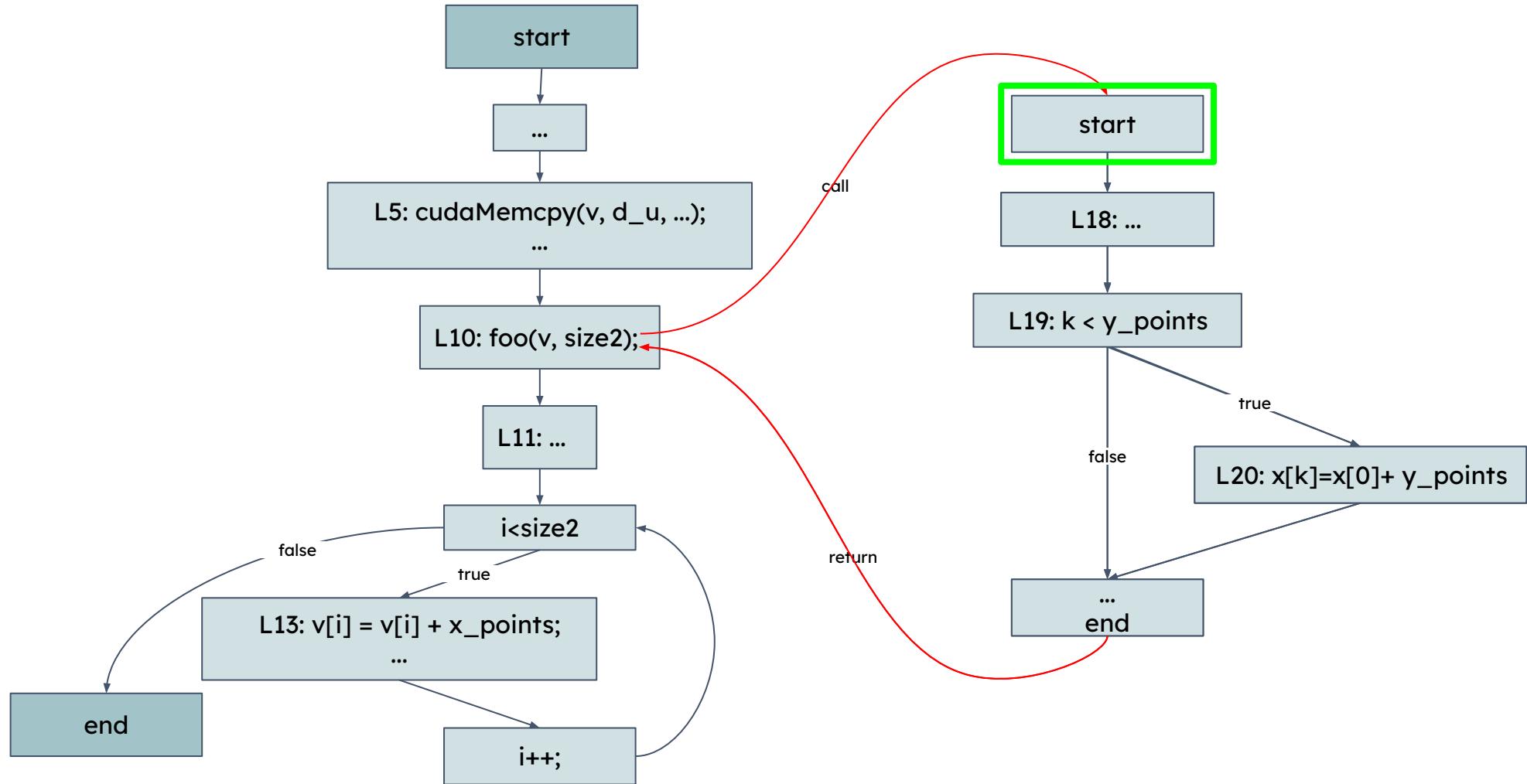
L5 $\mapsto \{?\}$



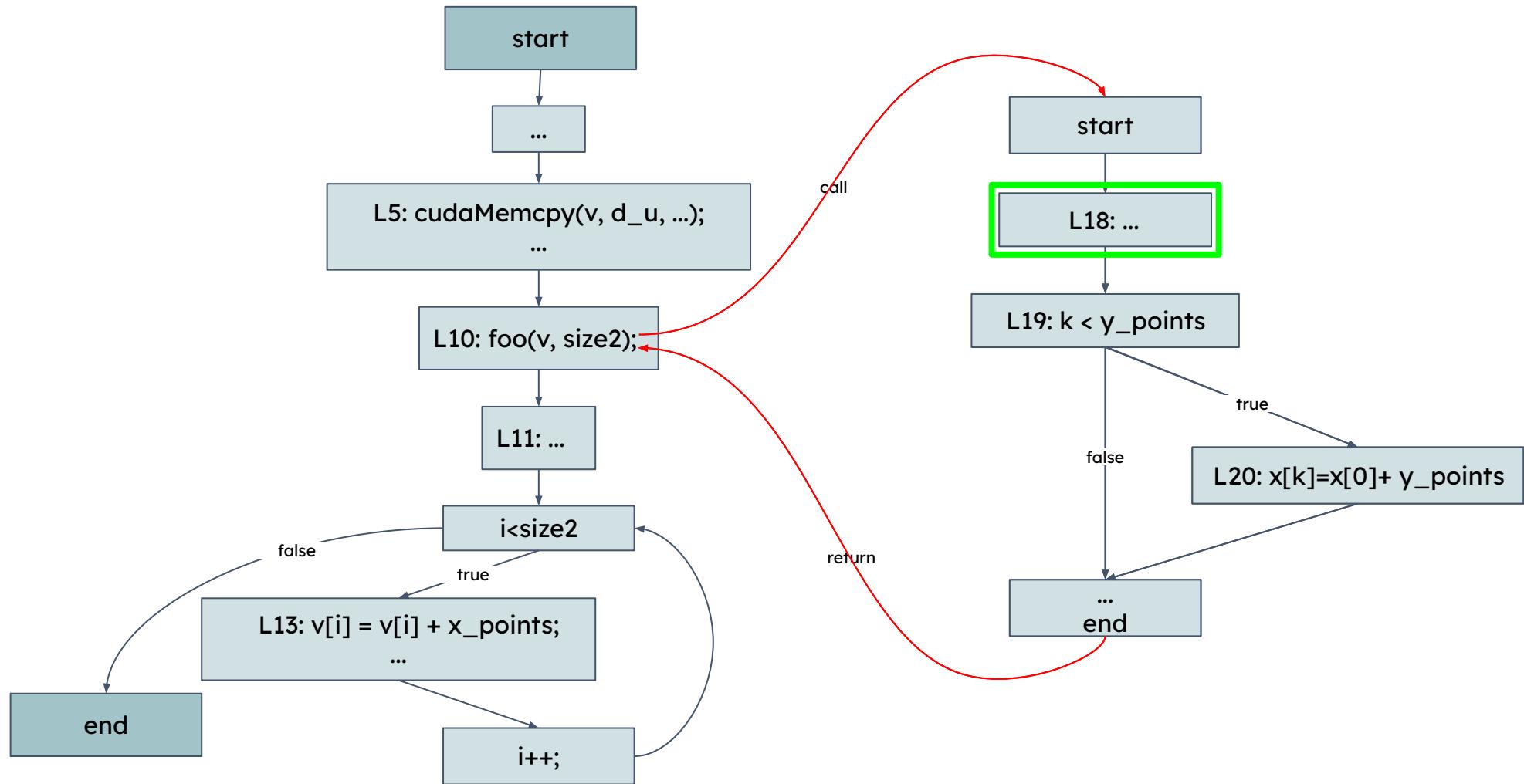
Step 1: Finding Target Locations (Contd.)



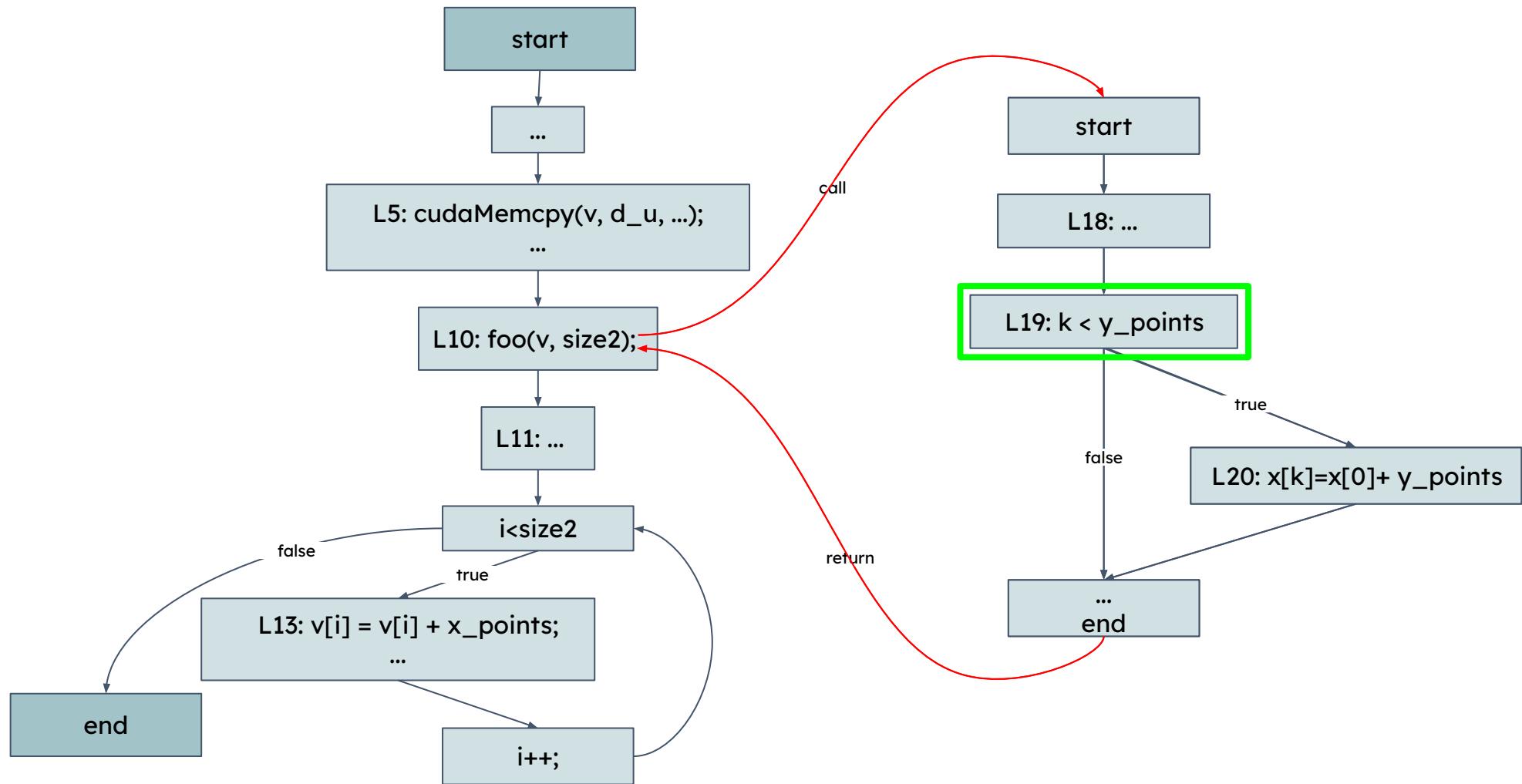
Step 1: Finding Target Locations (Contd.)



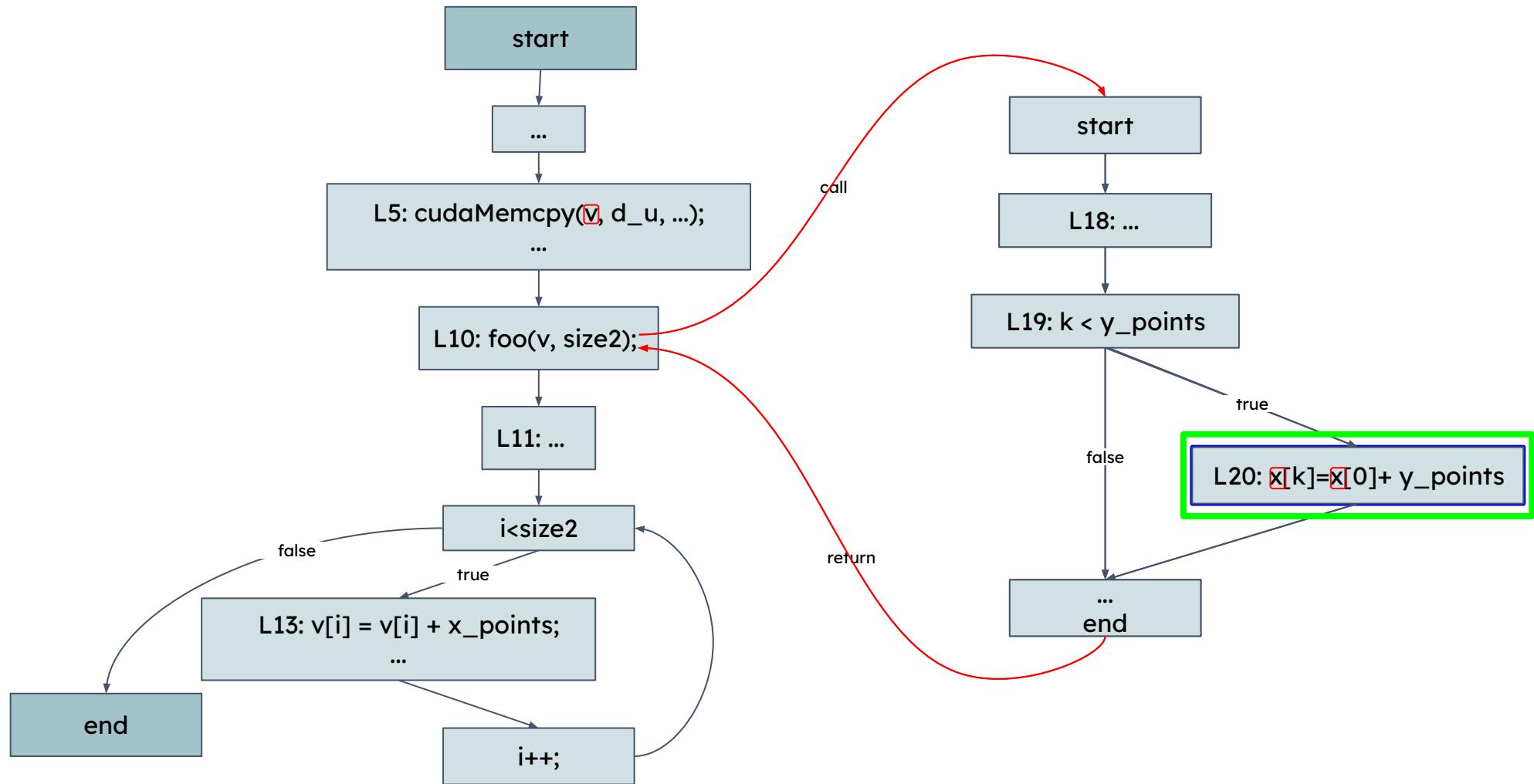
Step 1: Finding Target Locations (Contd.)



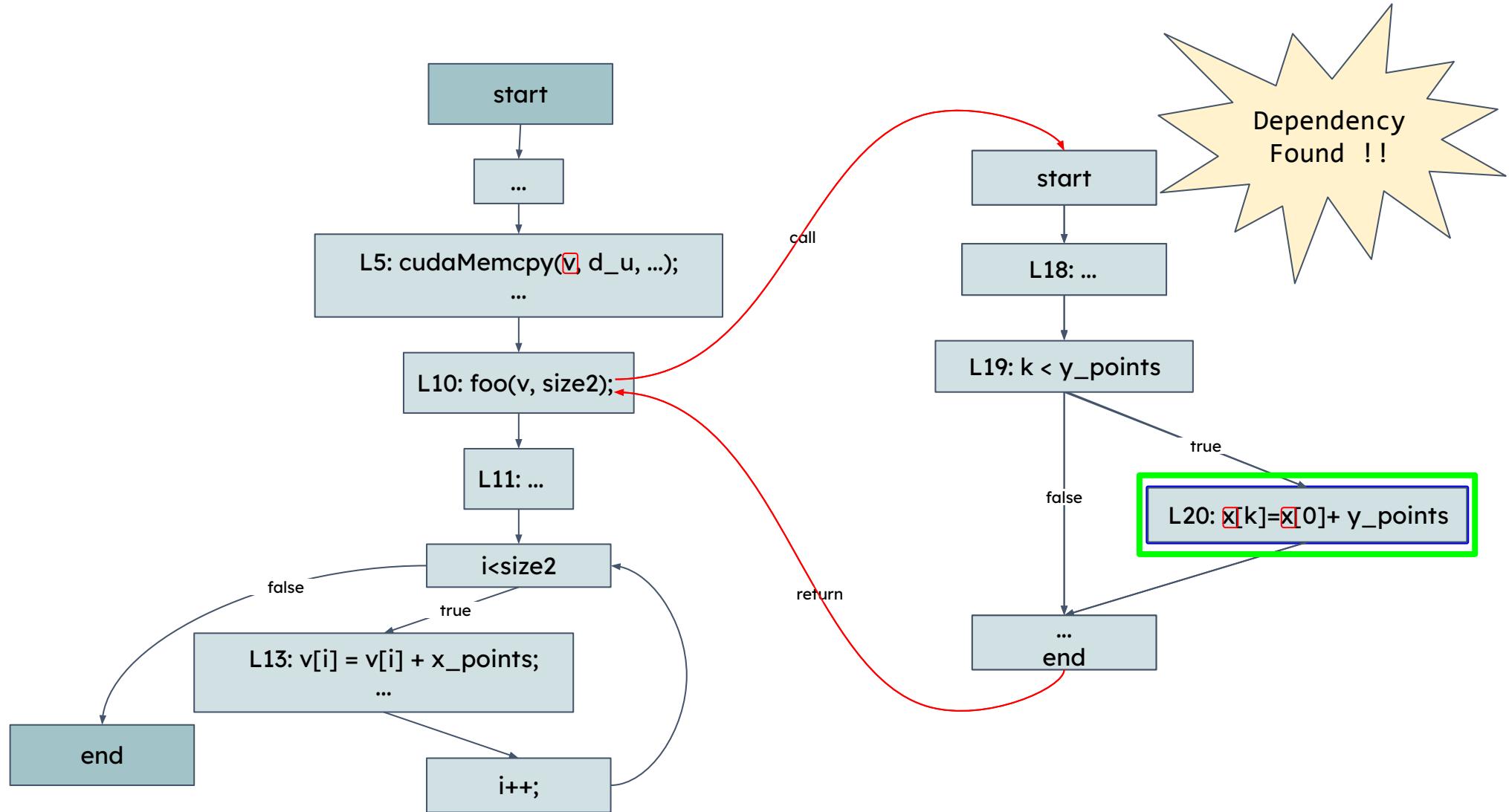
Step 1: Finding Target Locations (Contd.)



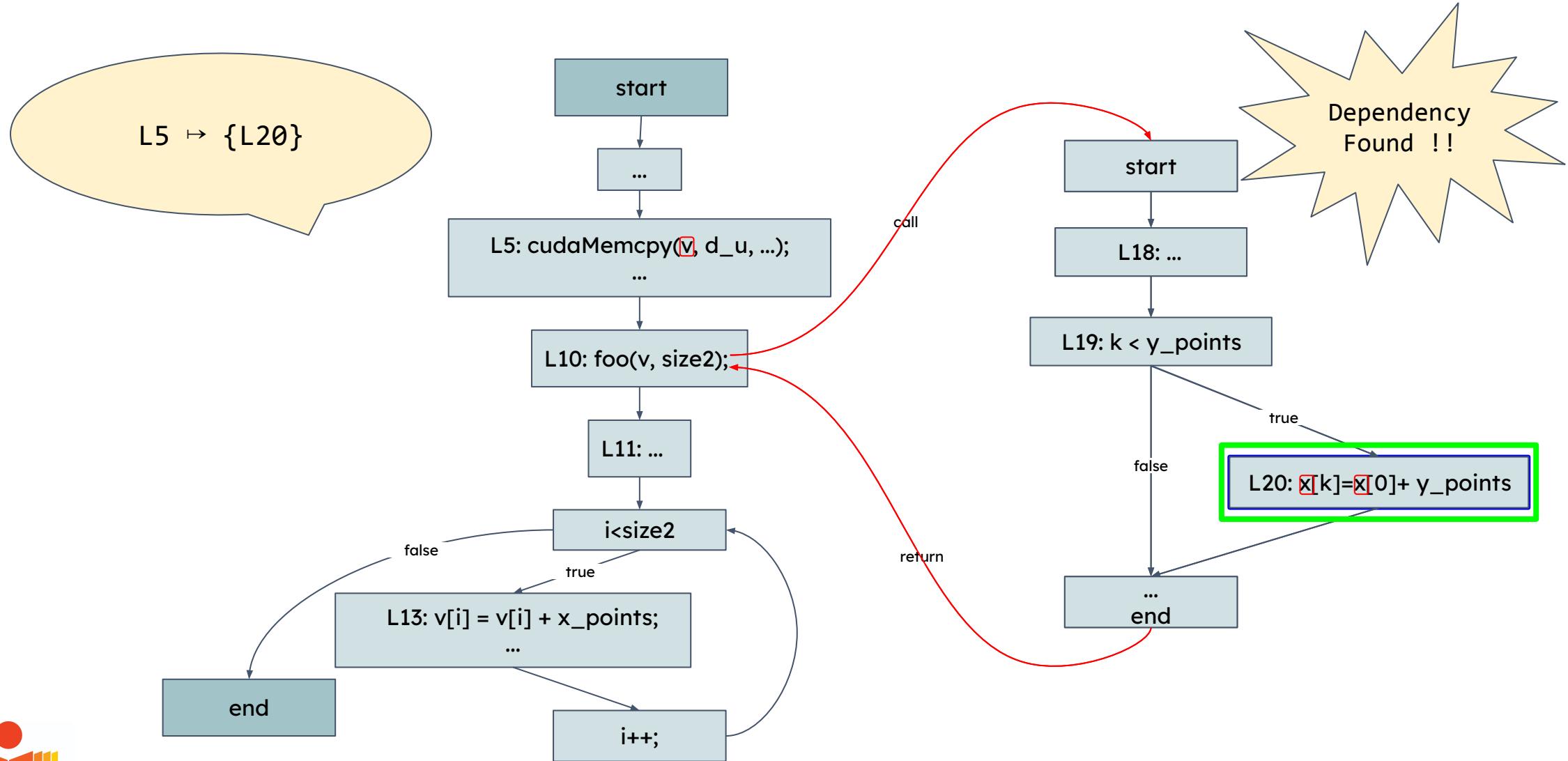
Step 1: Finding Target Locations (Contd.)



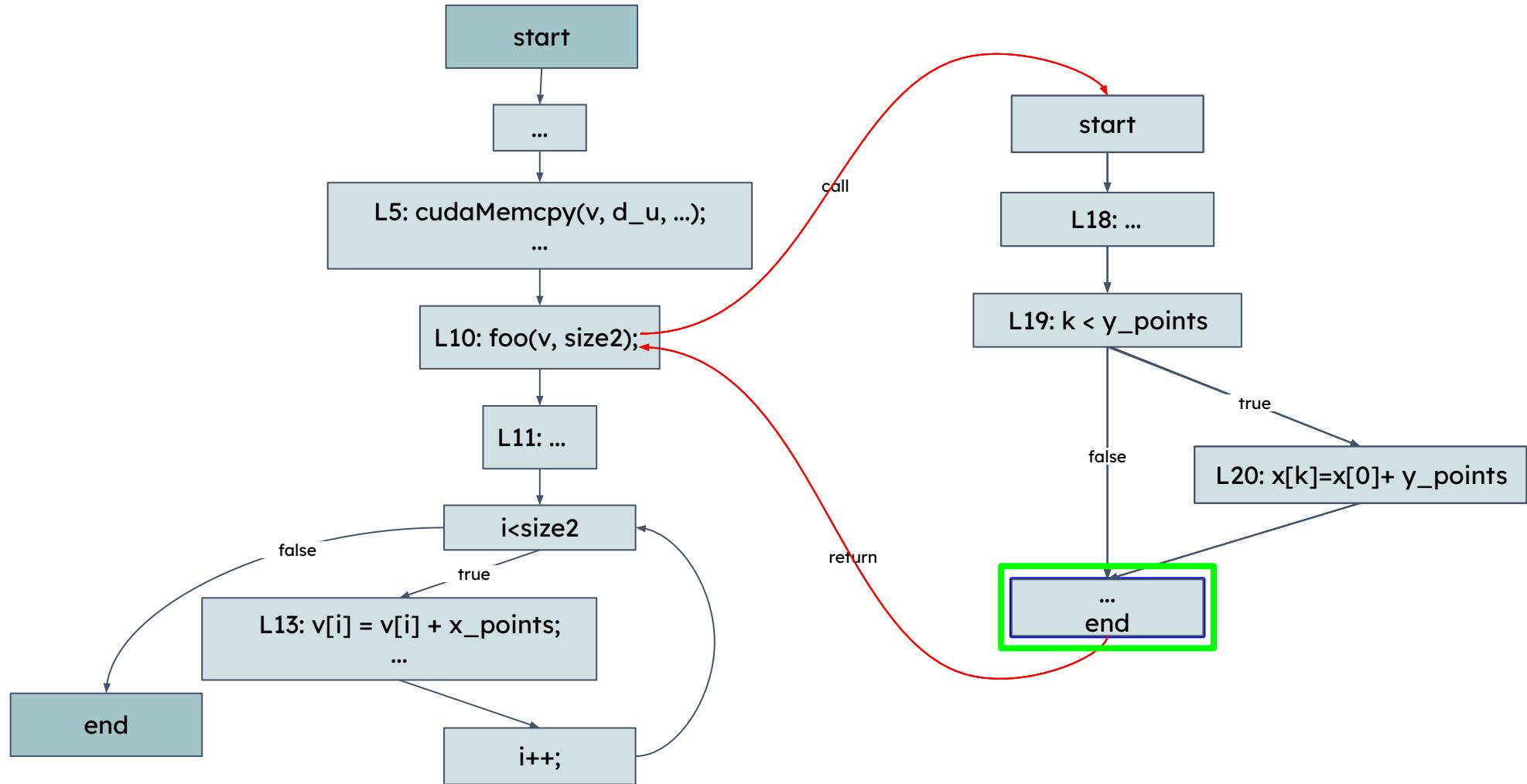
Step 1: Finding Target Locations (Contd.)



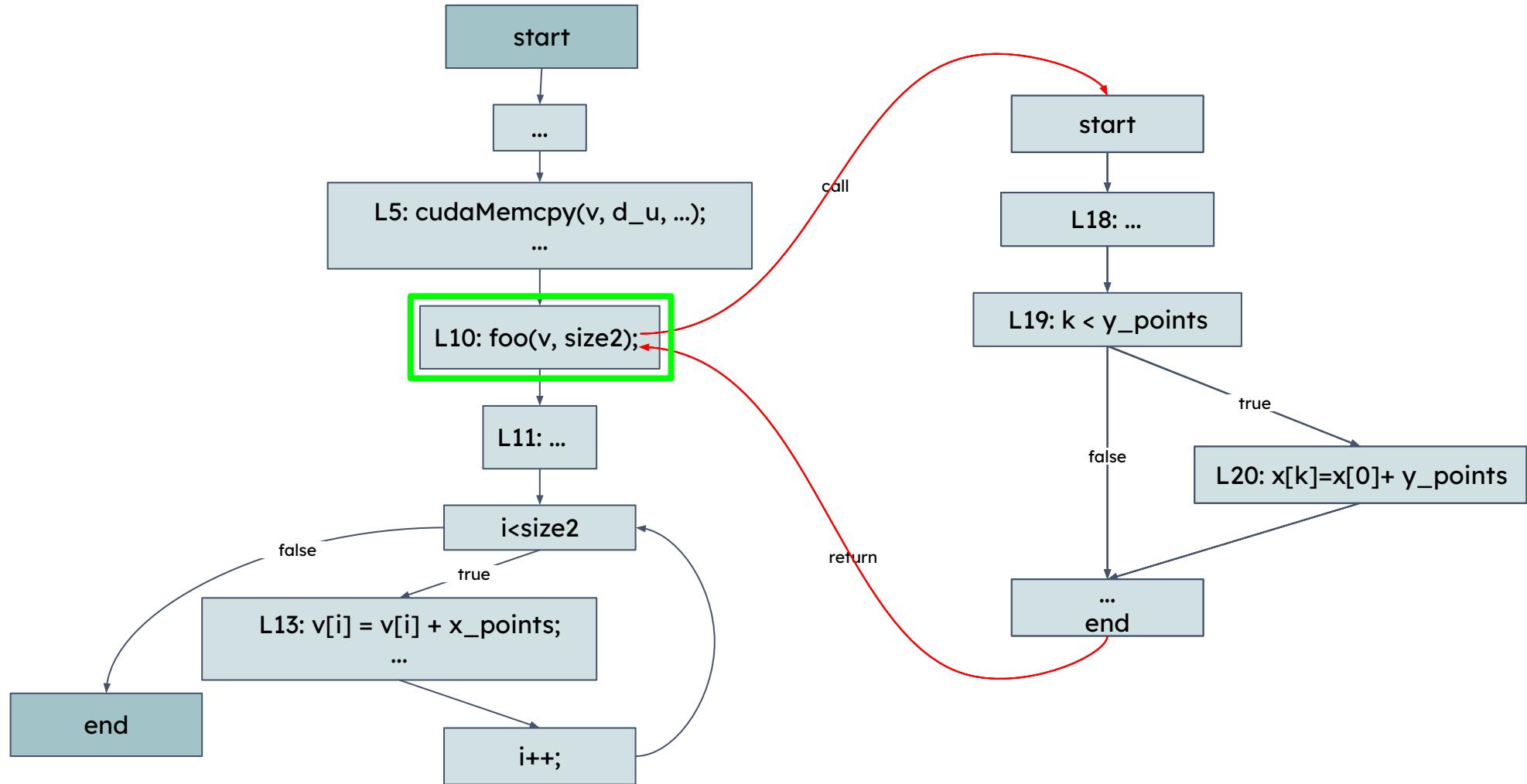
Step 1: Finding Target Locations (Contd.)



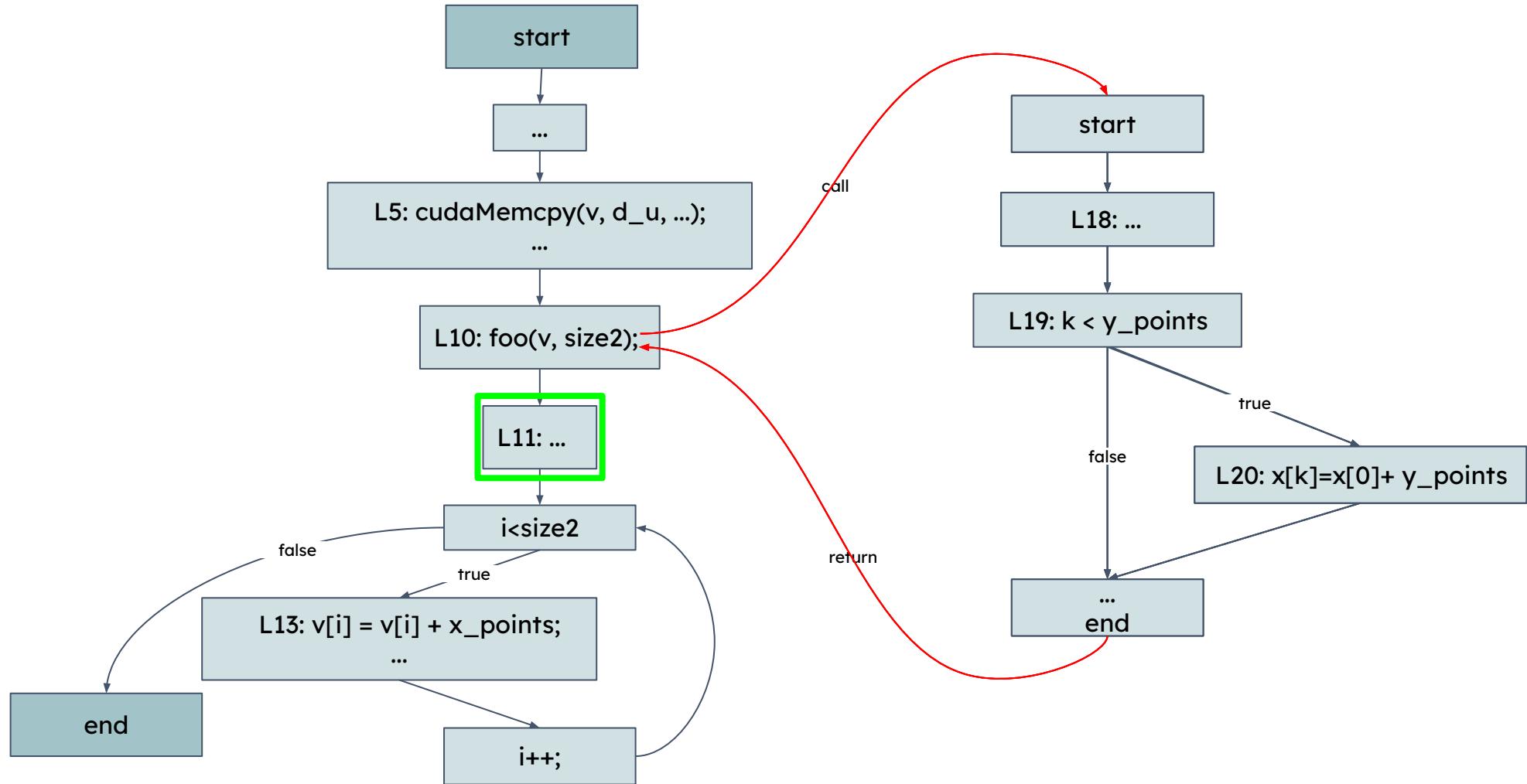
Step 1: Finding Target Locations (Contd.)



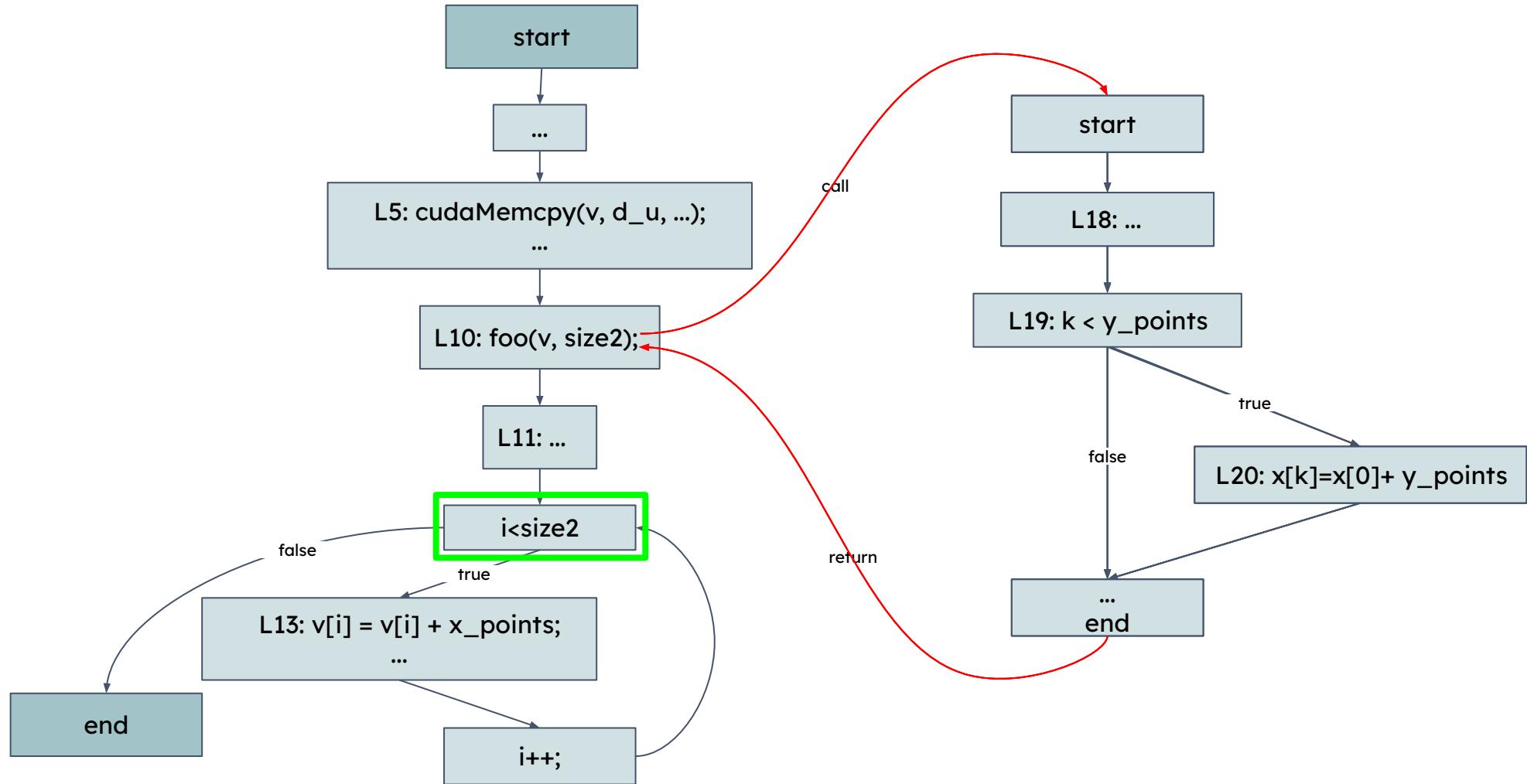
Step 1: Finding Target Locations (Contd.)



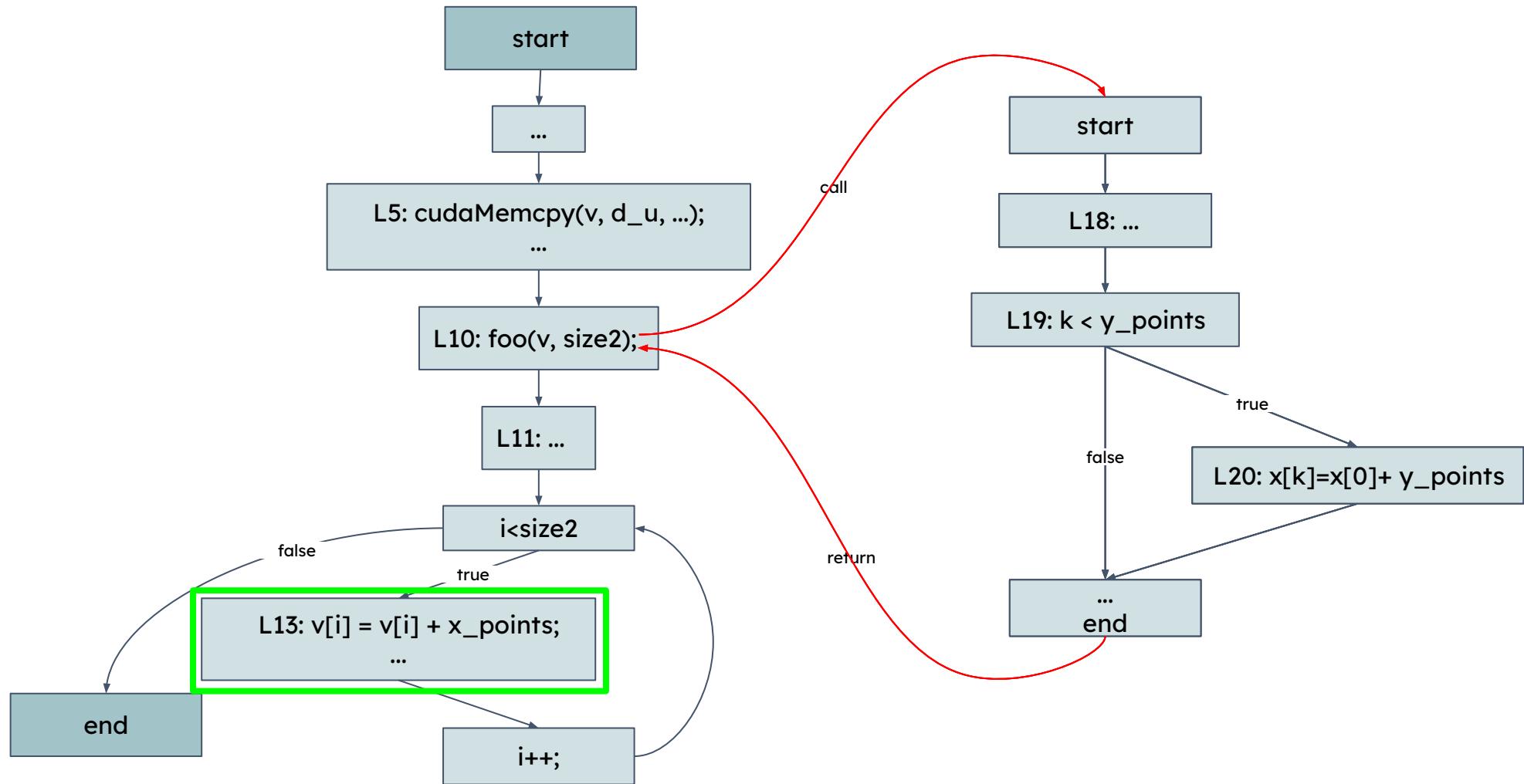
Step 1: Finding Target Locations (Contd.)



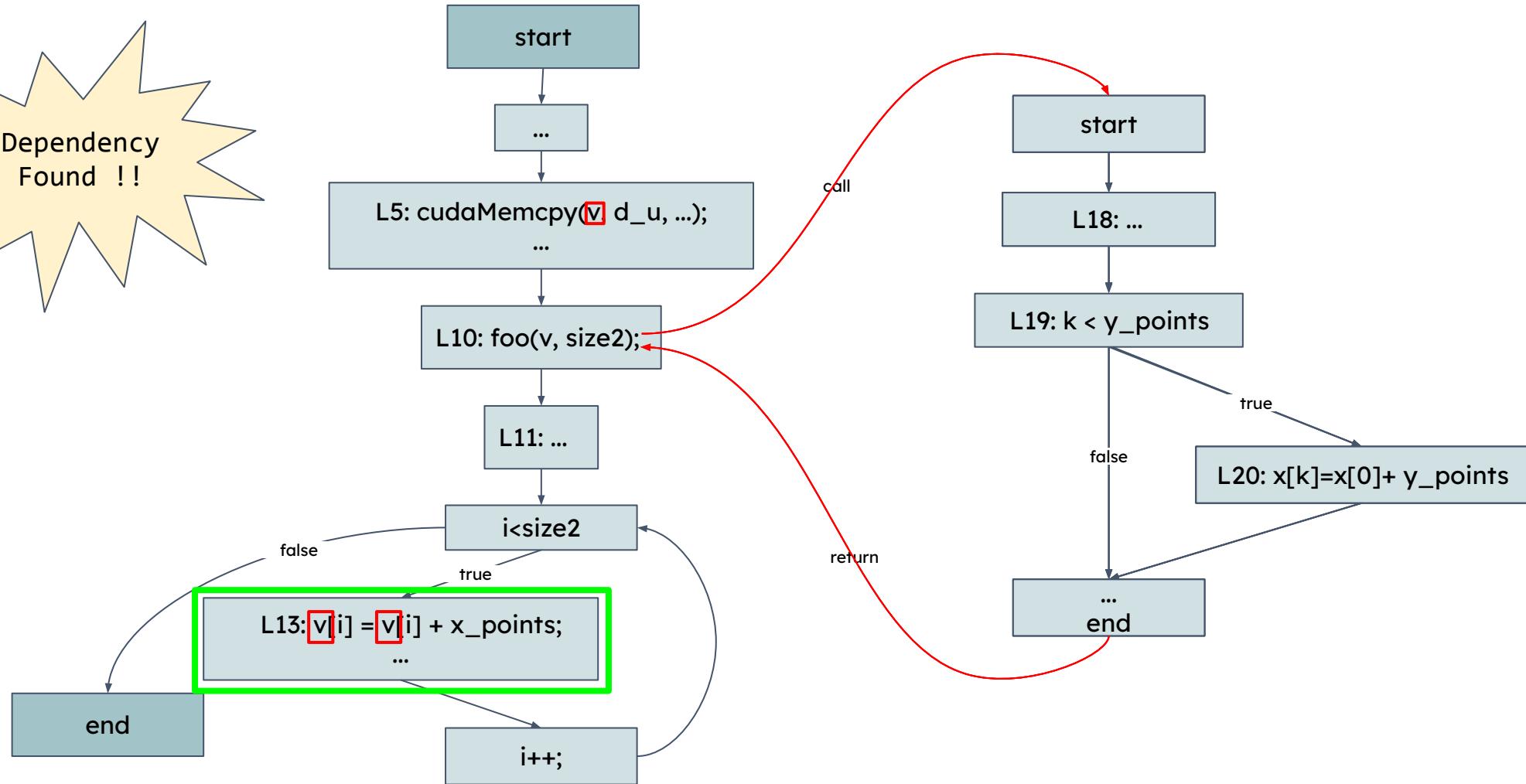
Step 1: Finding Target Locations (Contd.)



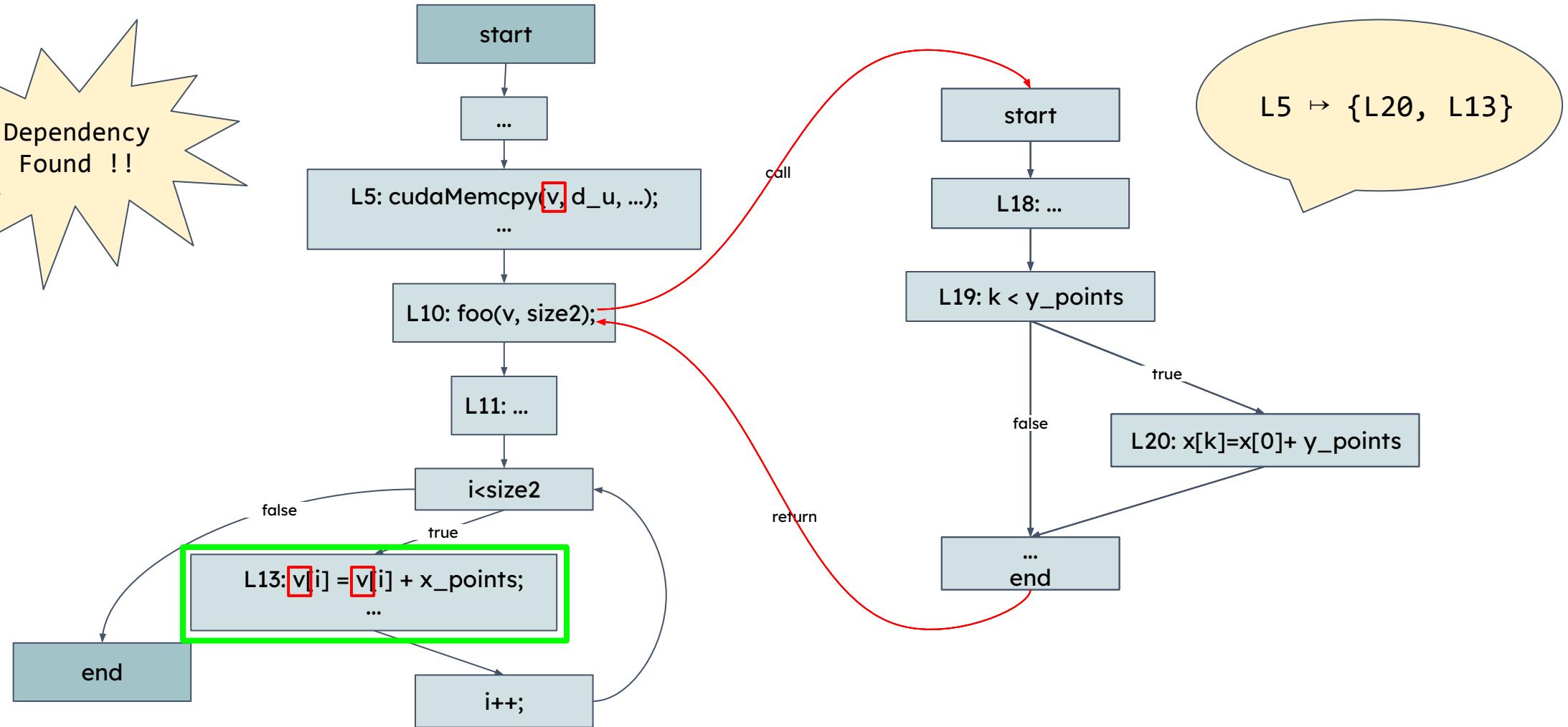
Step 1: Finding Target Locations (Contd.)



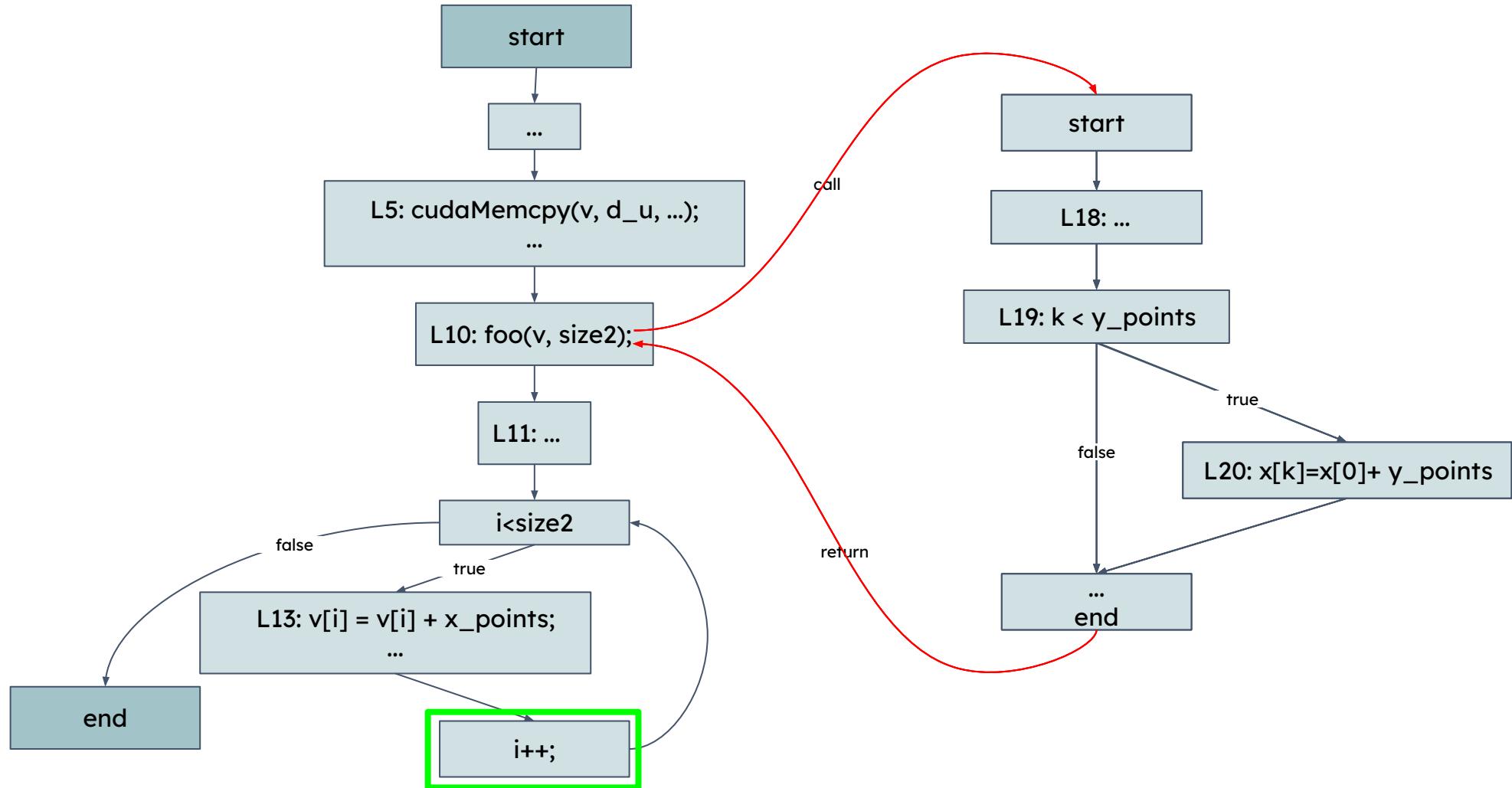
Step 1: Finding Target Locations (Contd.)



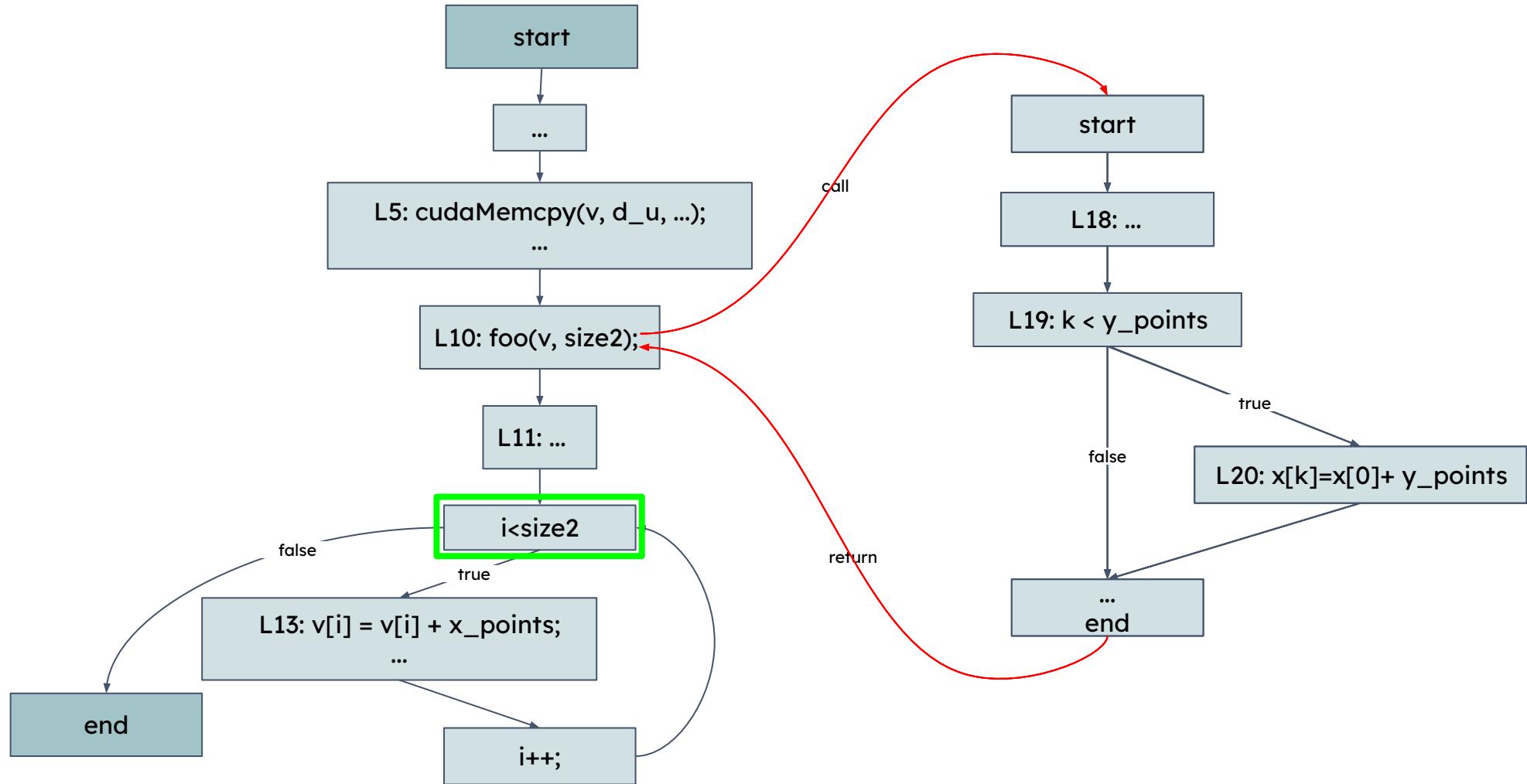
Step 1: Finding Target Locations (Contd.)



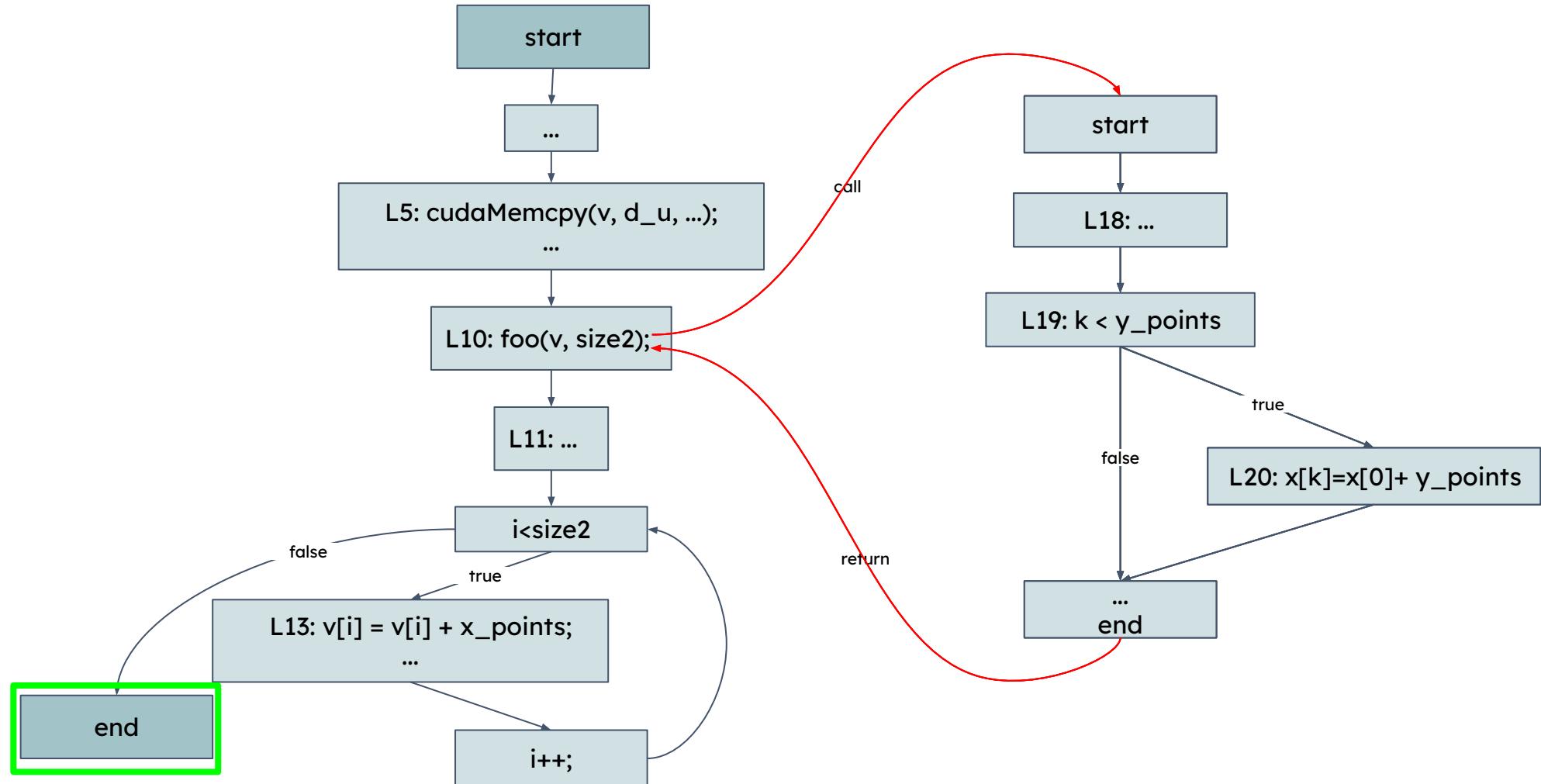
Step 1: Finding Target Locations (Contd.)



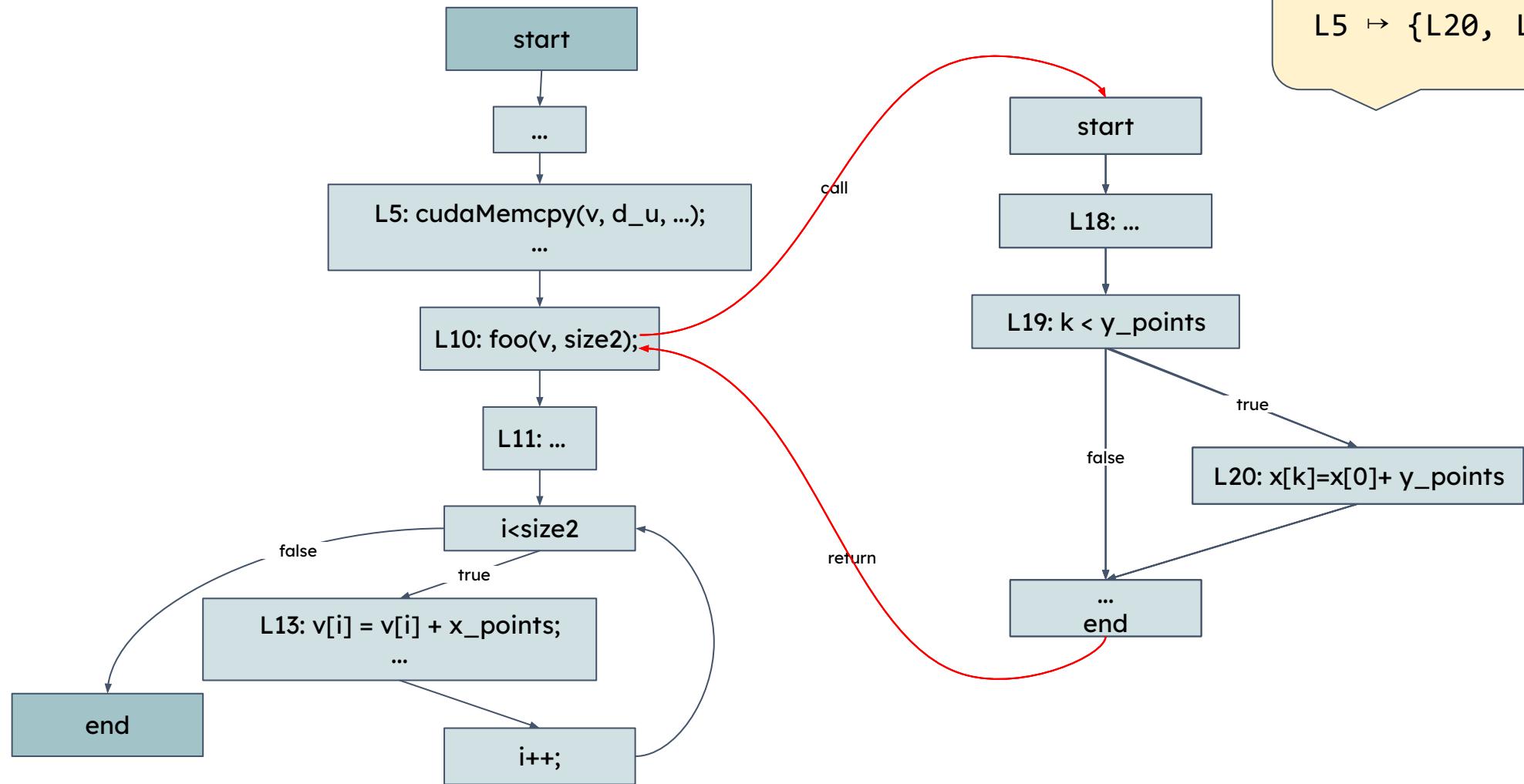
Step 1: Finding Target Locations (Contd.)



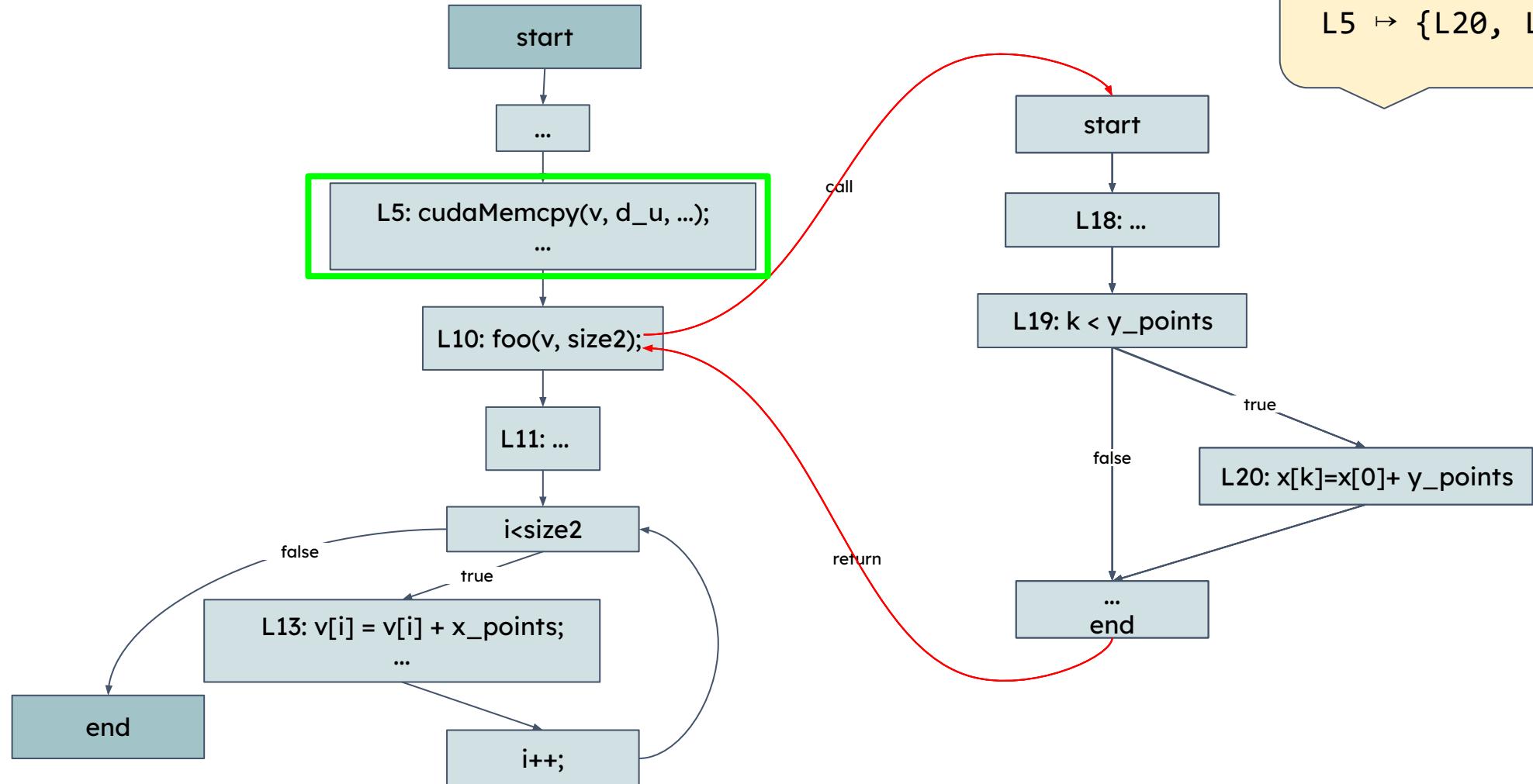
Step 1: Finding Target Locations (Contd.)



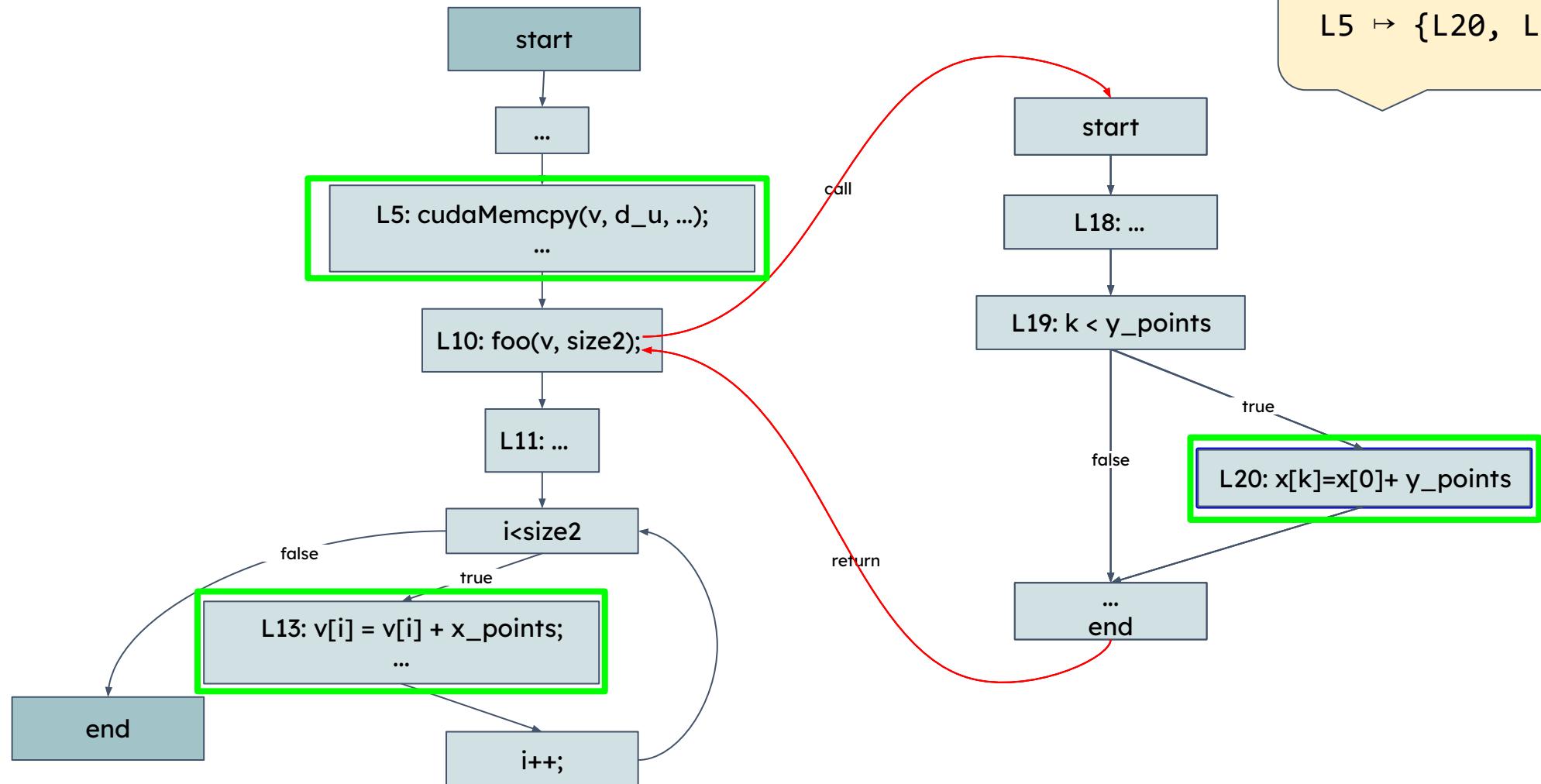
Step 2: Updating Target Locations: Motivation



Step 2: Updating Target Locations: Motivation

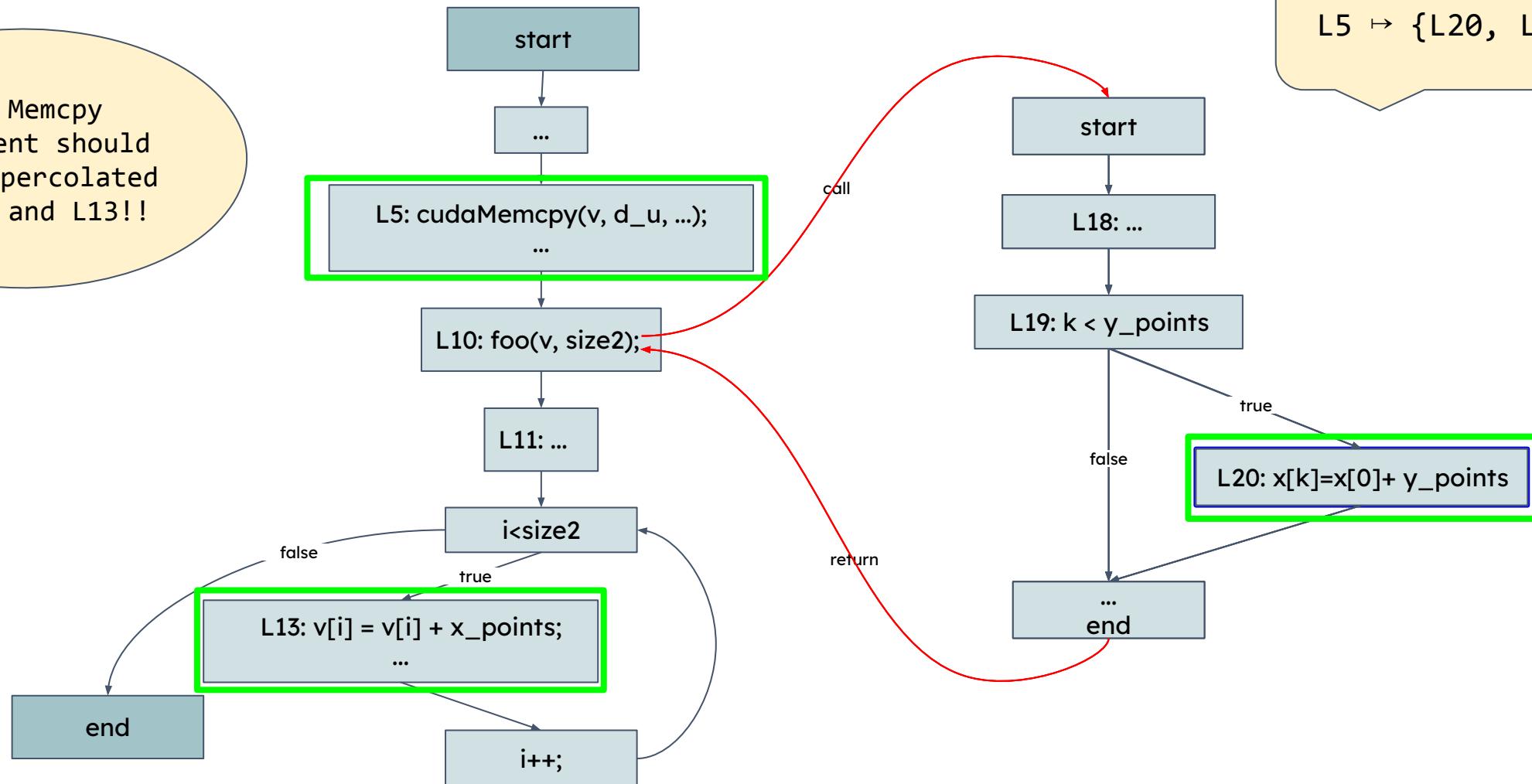


Step 2: Updating Target Locations: Motivation

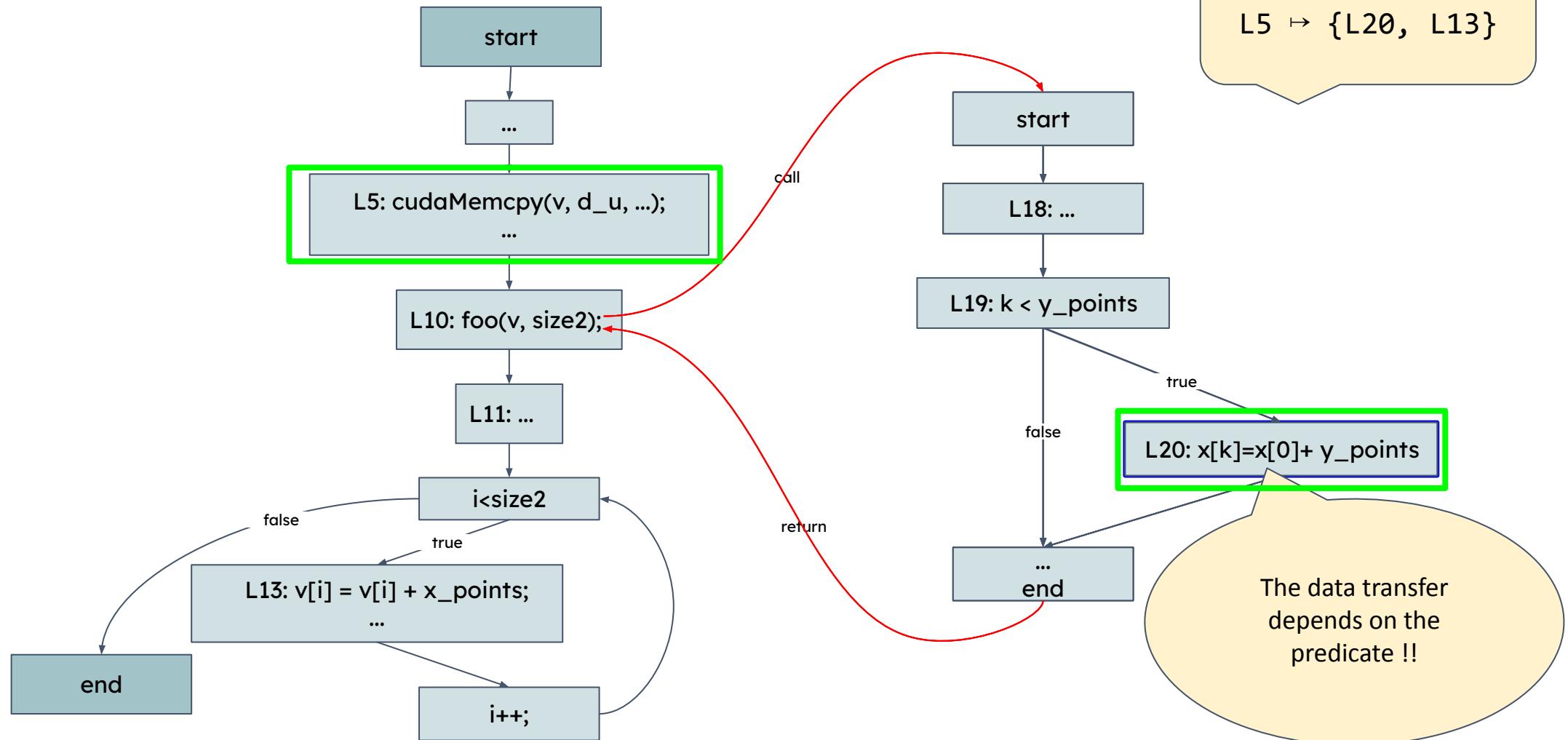


Step 2: Updating Target Locations: Motivation

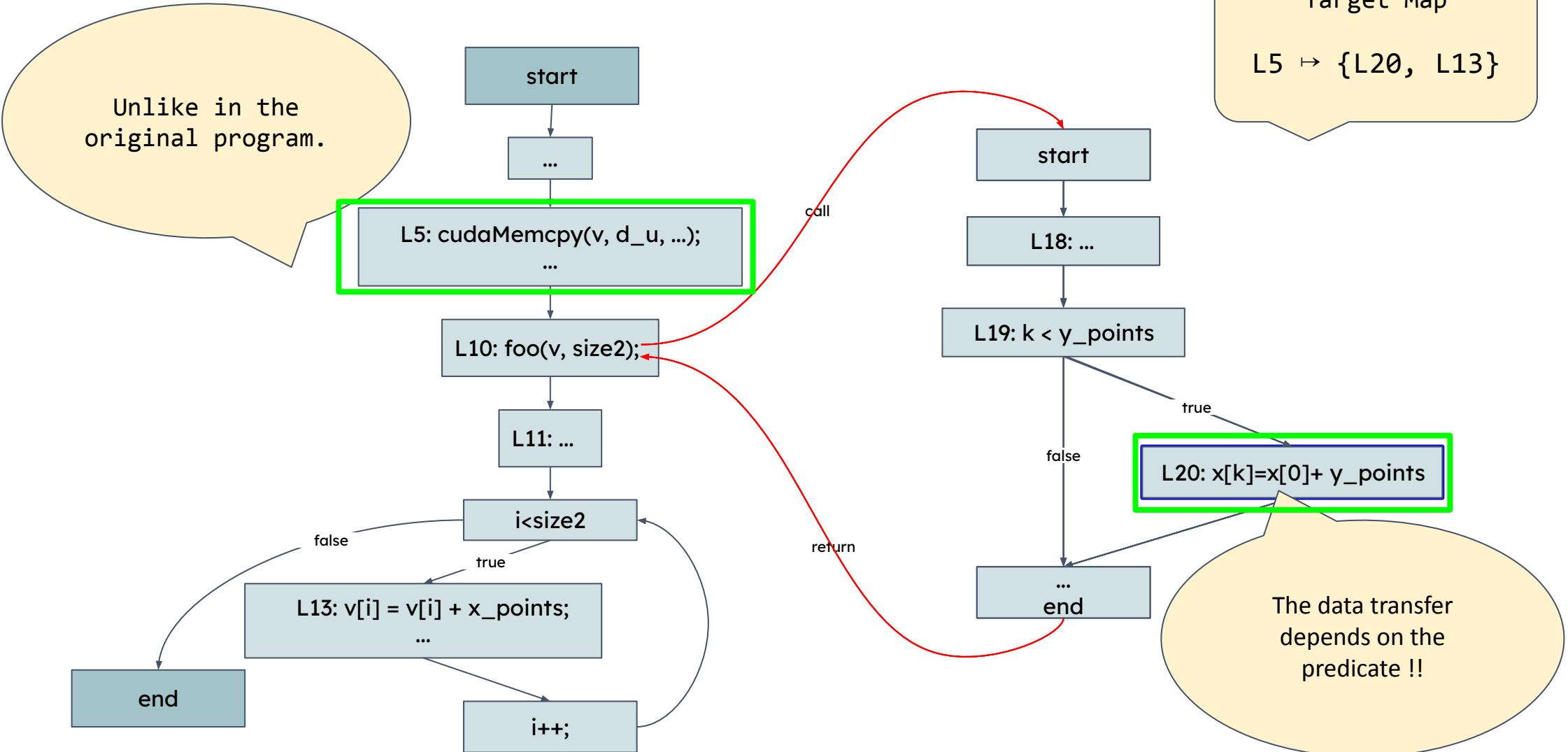
The Memcpy statement should not be percolated to L20 and L13!!



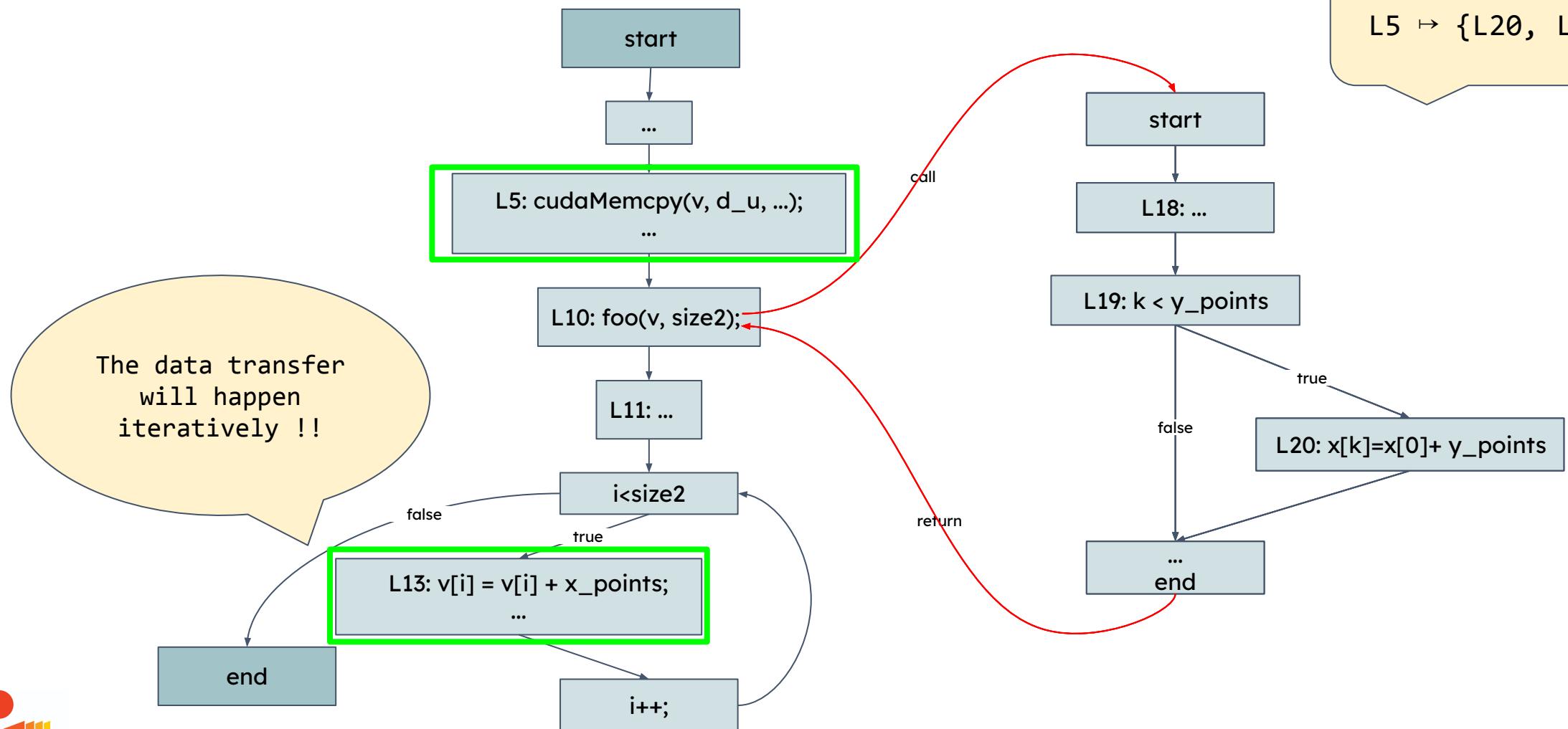
Step 2: Updating Target Locations: Motivation



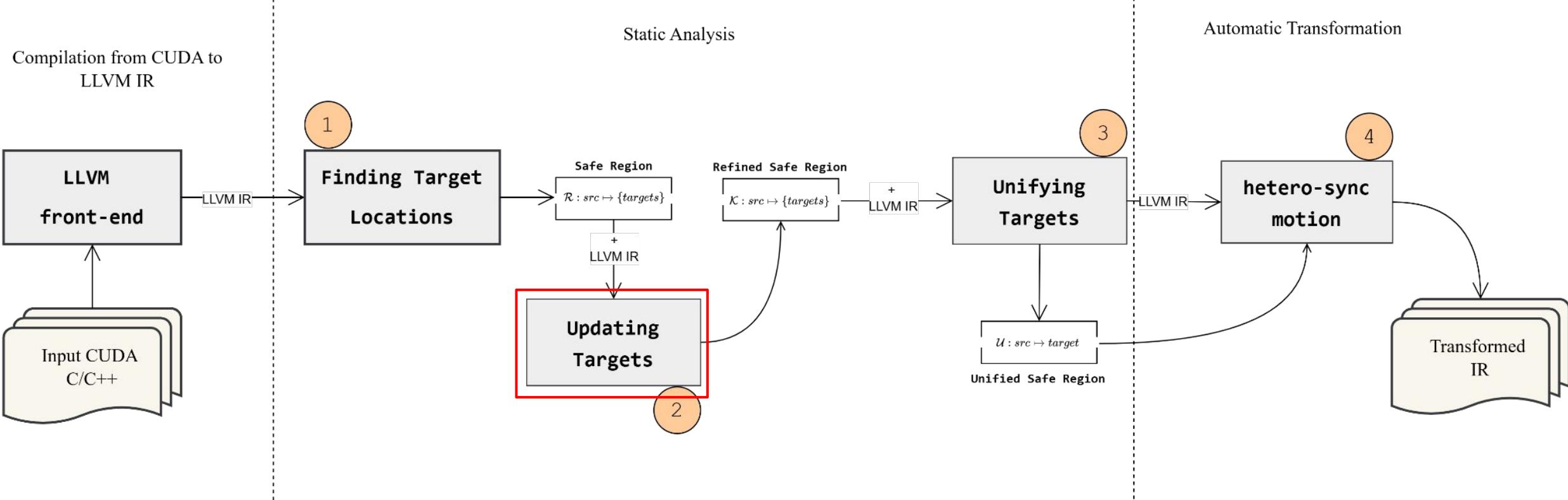
Step 2: Updating Target Locations: Motivation



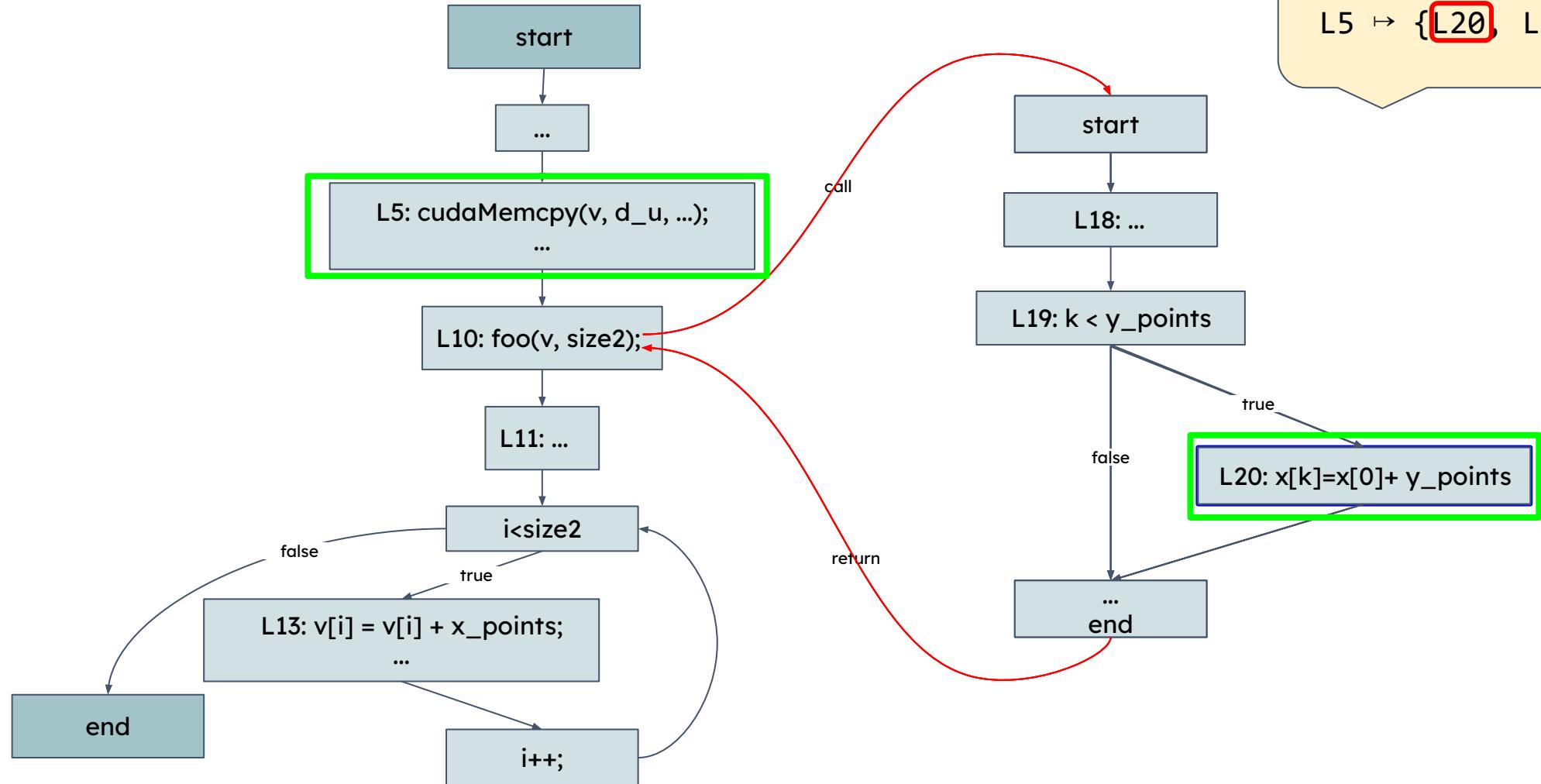
Step 2: Updating Target Locations: Motivation



Step 2: Updating Target Locations

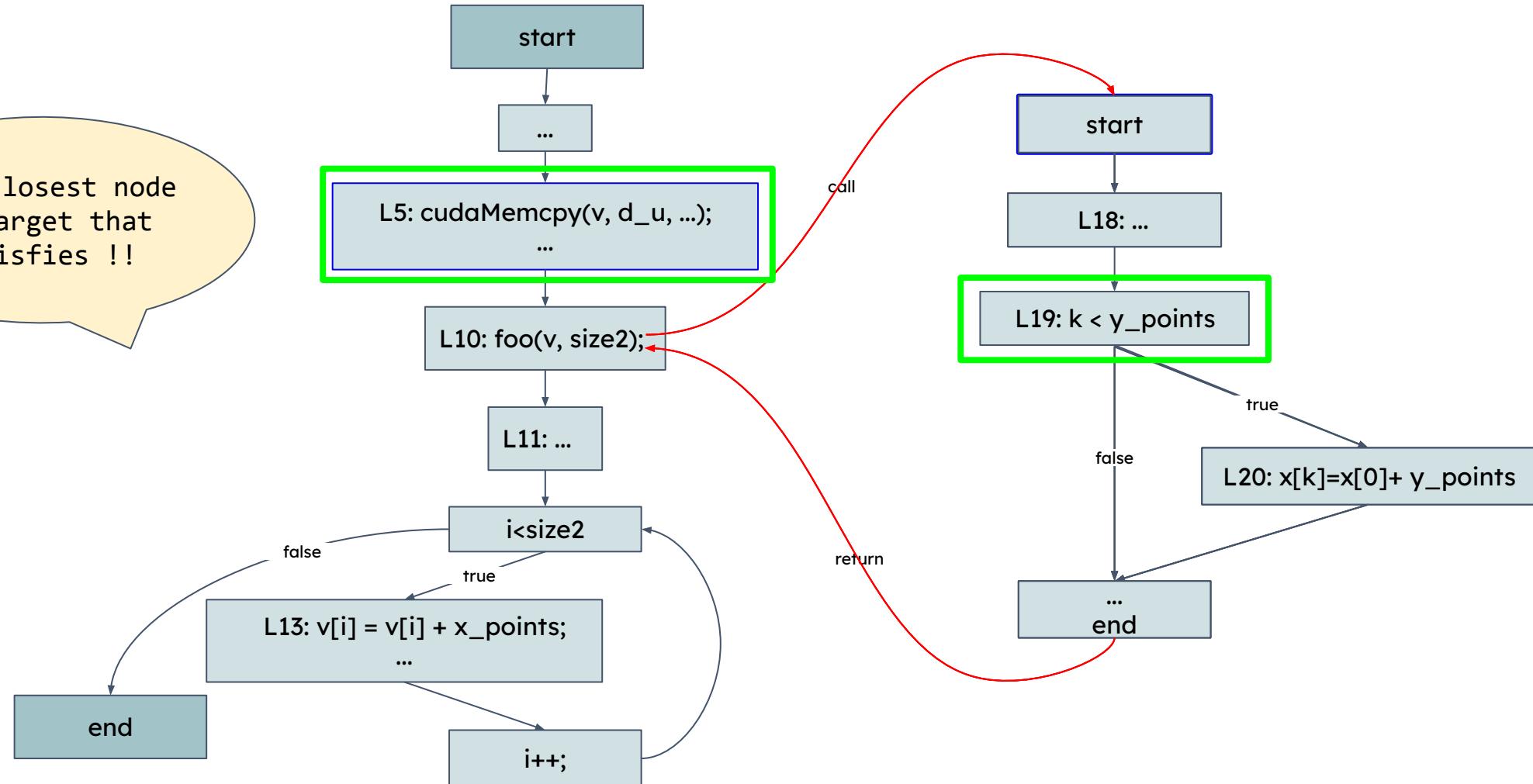


Step 2: Updating Target Locations: Motivation



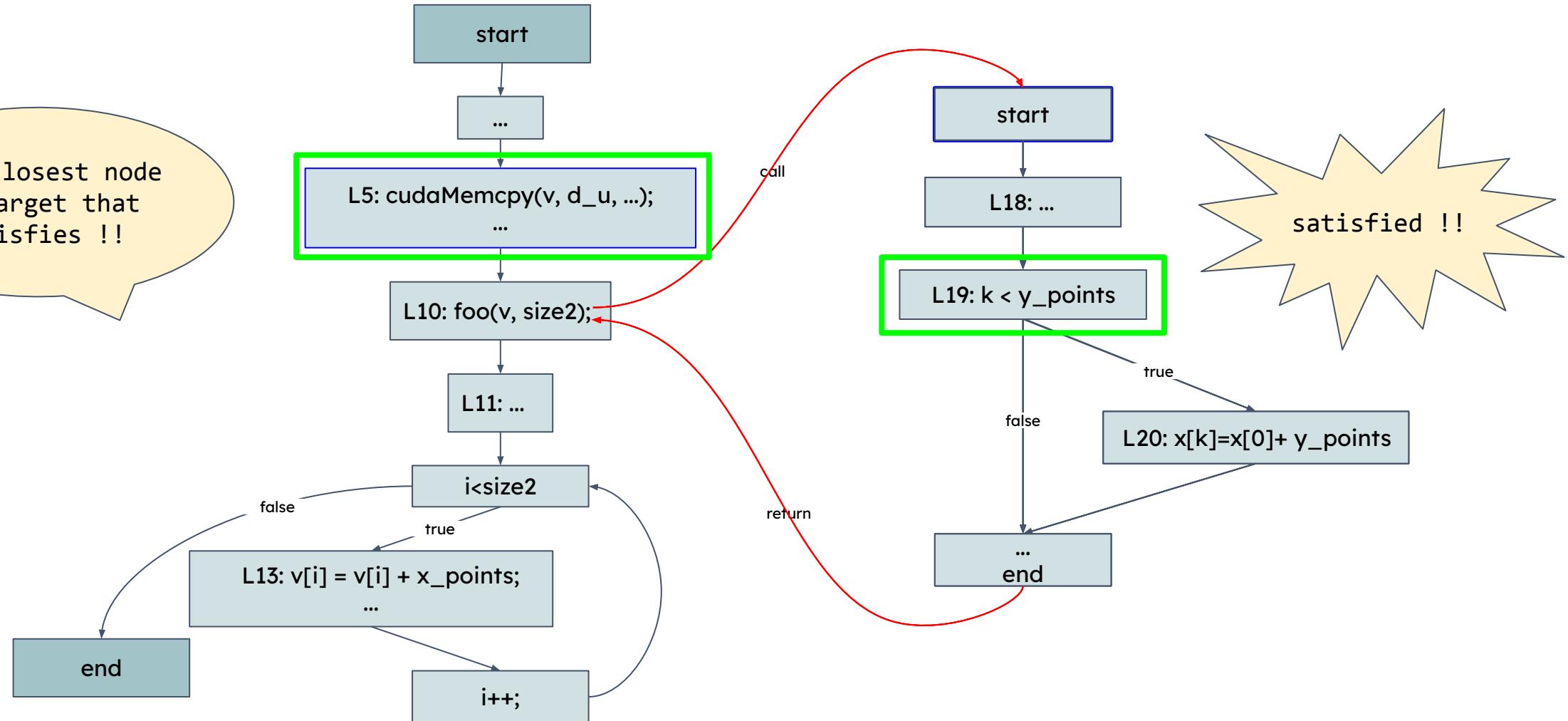
Step 2: Updating Target Locations (Contd.)

Find closest node
to target that
satisfies !!



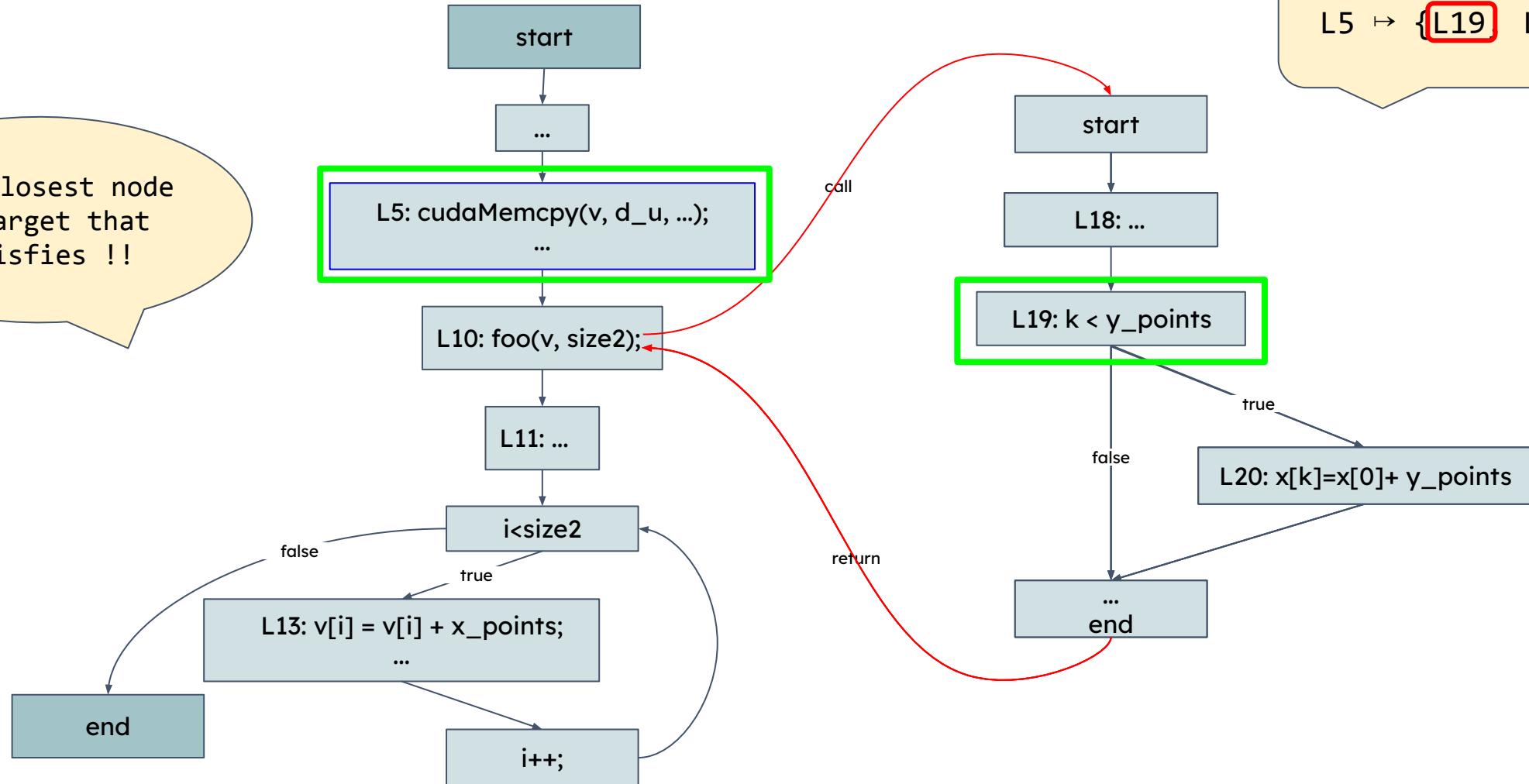
Step 2: Updating Target Locations (Contd.)

Find closest node
to target that
satisfies !!



Step 2: Updating Target Locations (Contd.)

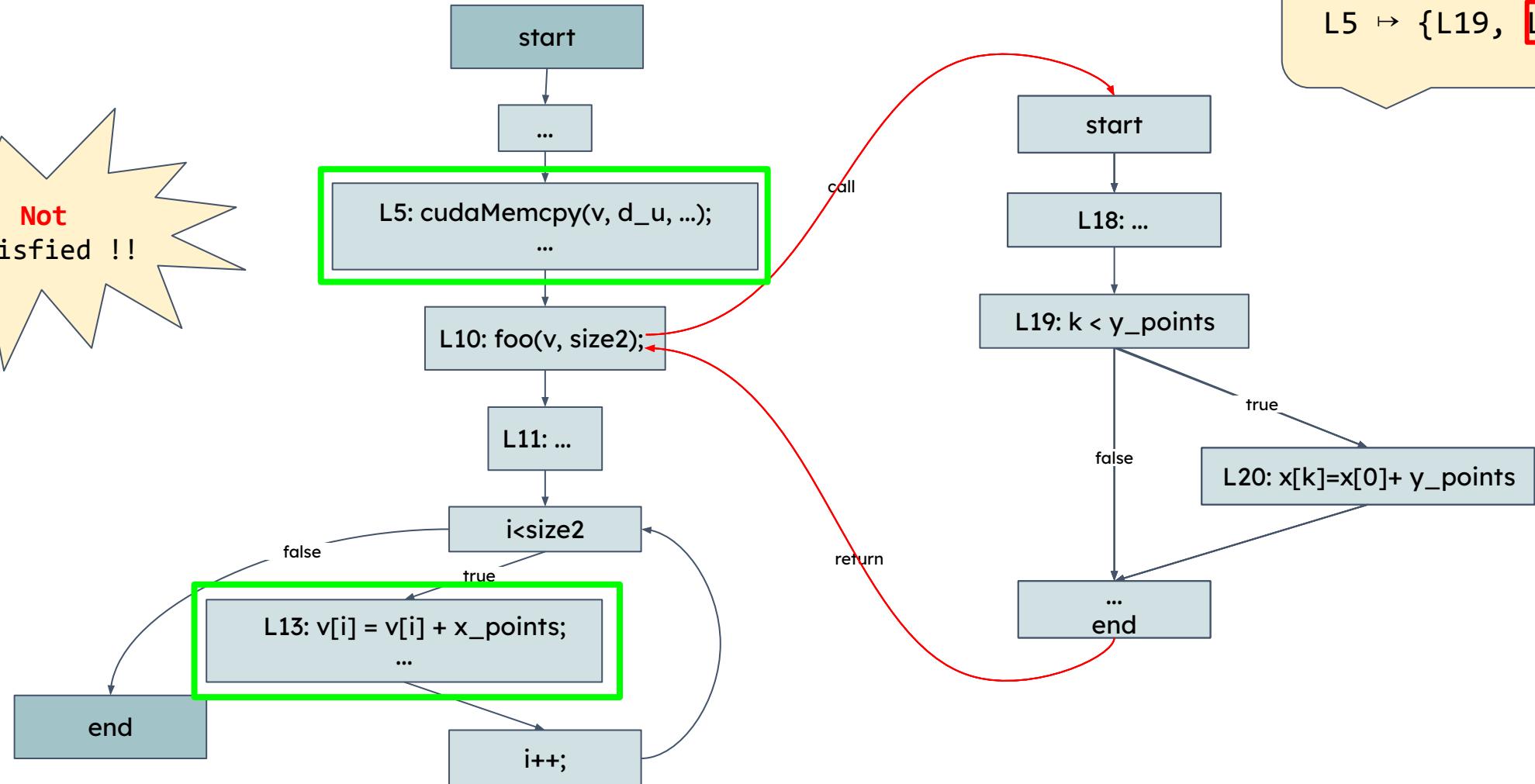
Find closest node
to target that
satisfies !!



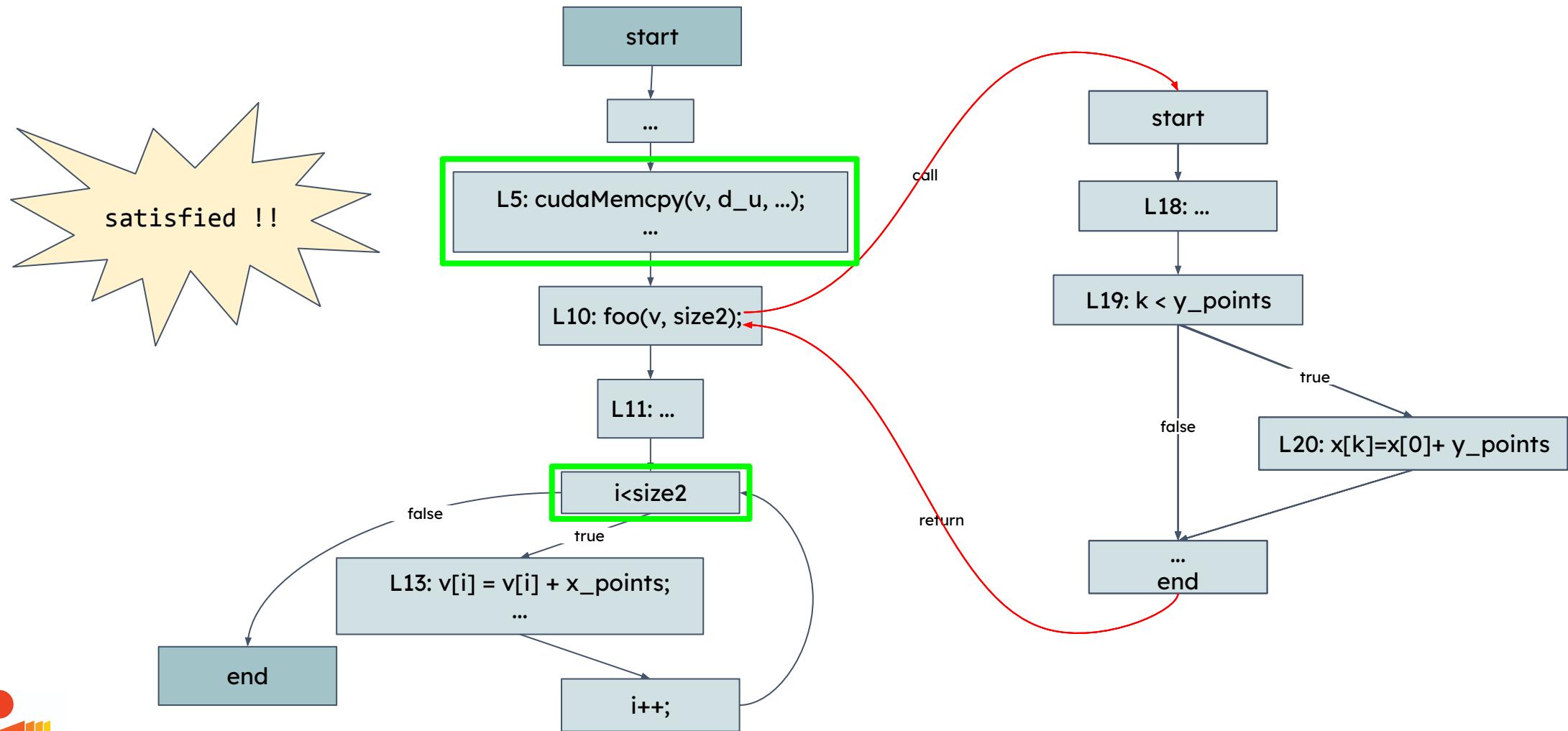
Target Map

$L5 \mapsto \{L19, L13\}$

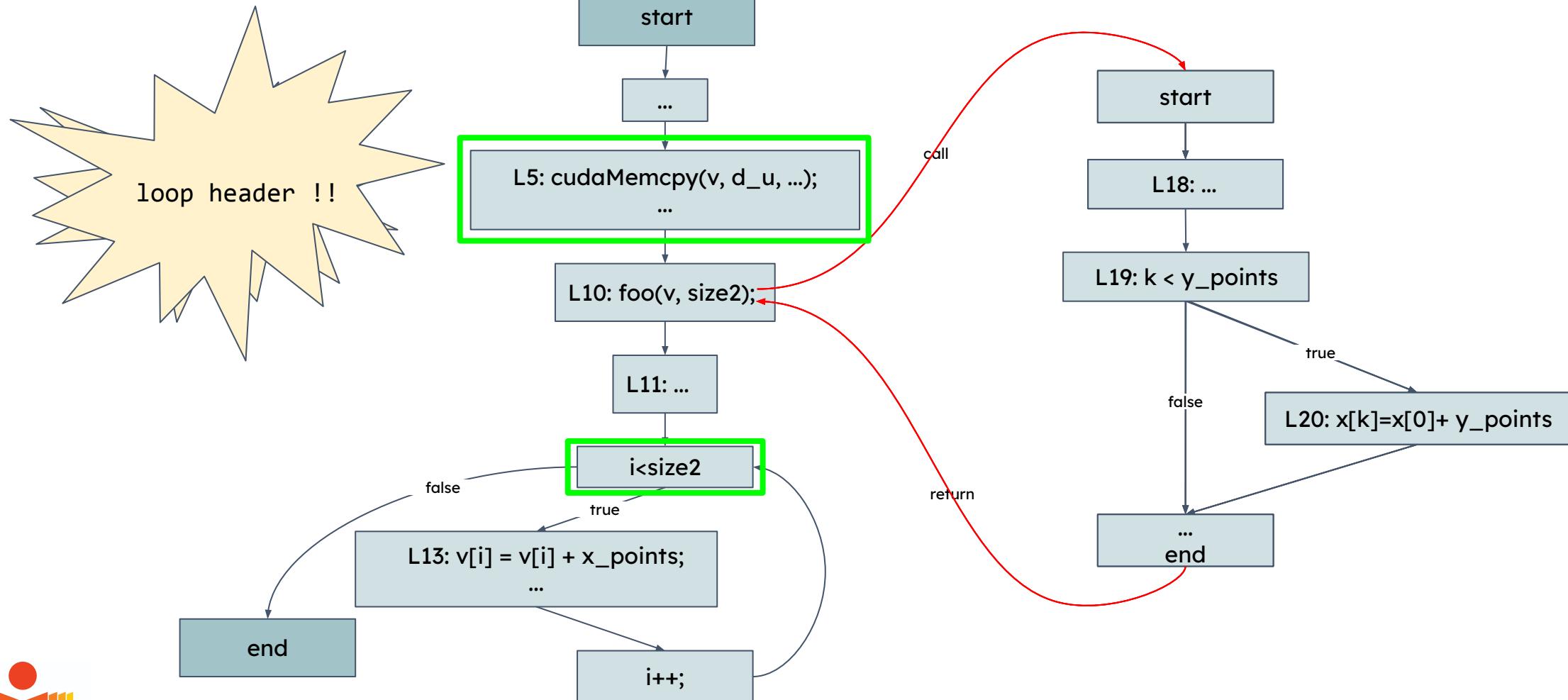
Step 2: Updating Target Locations (Contd.)



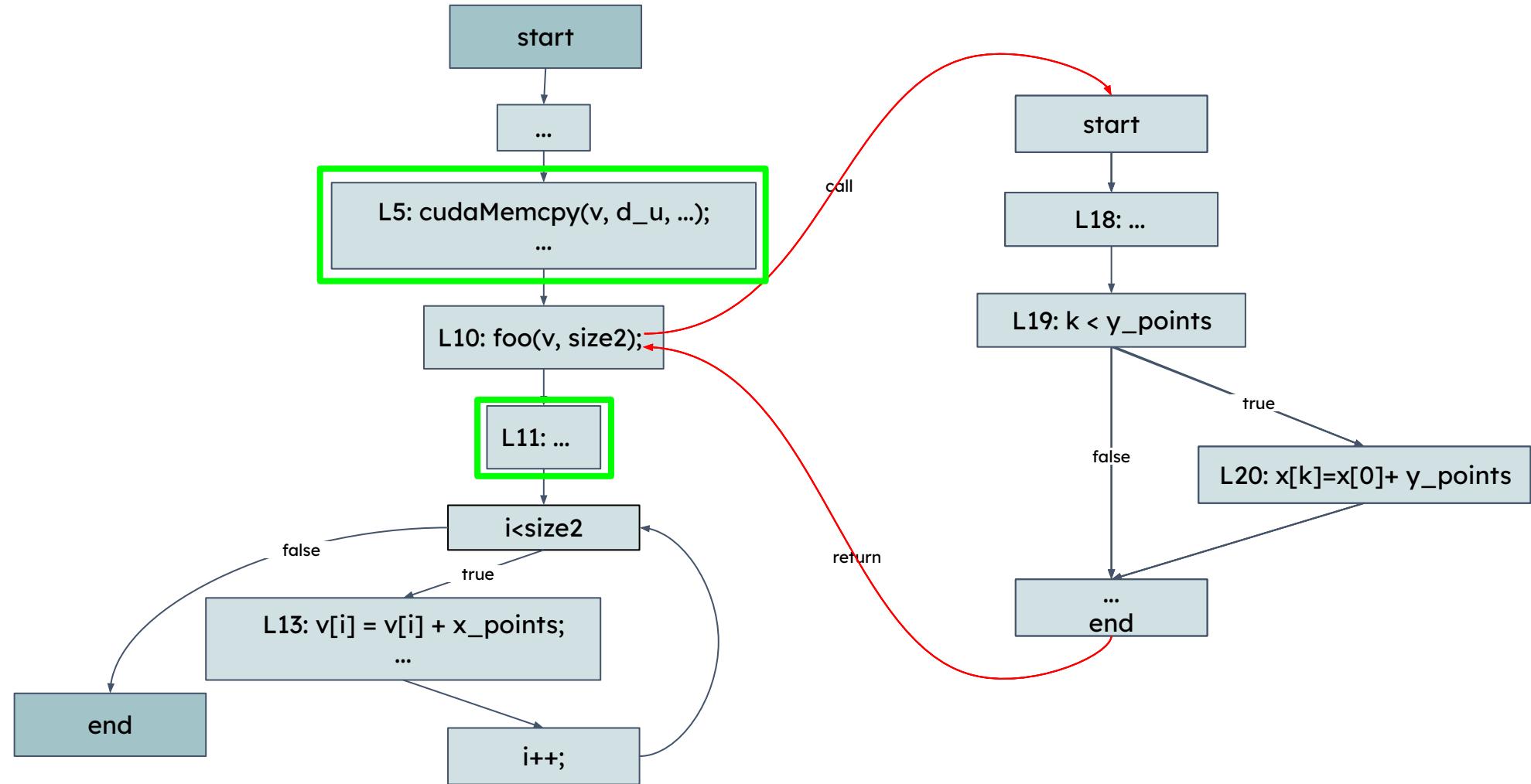
Step 2: Updating Target Locations (Contd.)



Step 2: Updating Target Locations (Contd.)

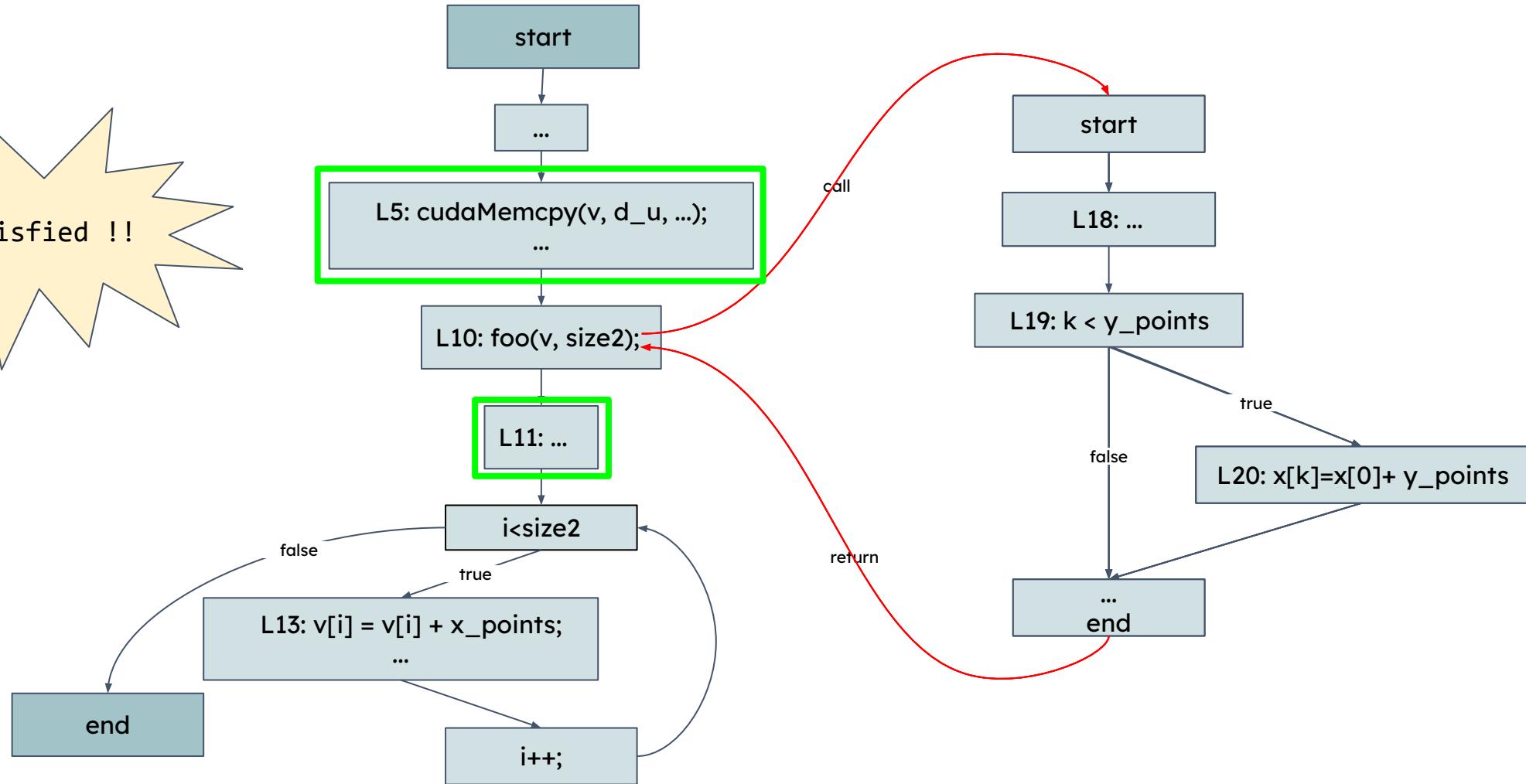


Step 2: Updating Target Locations (Contd.)

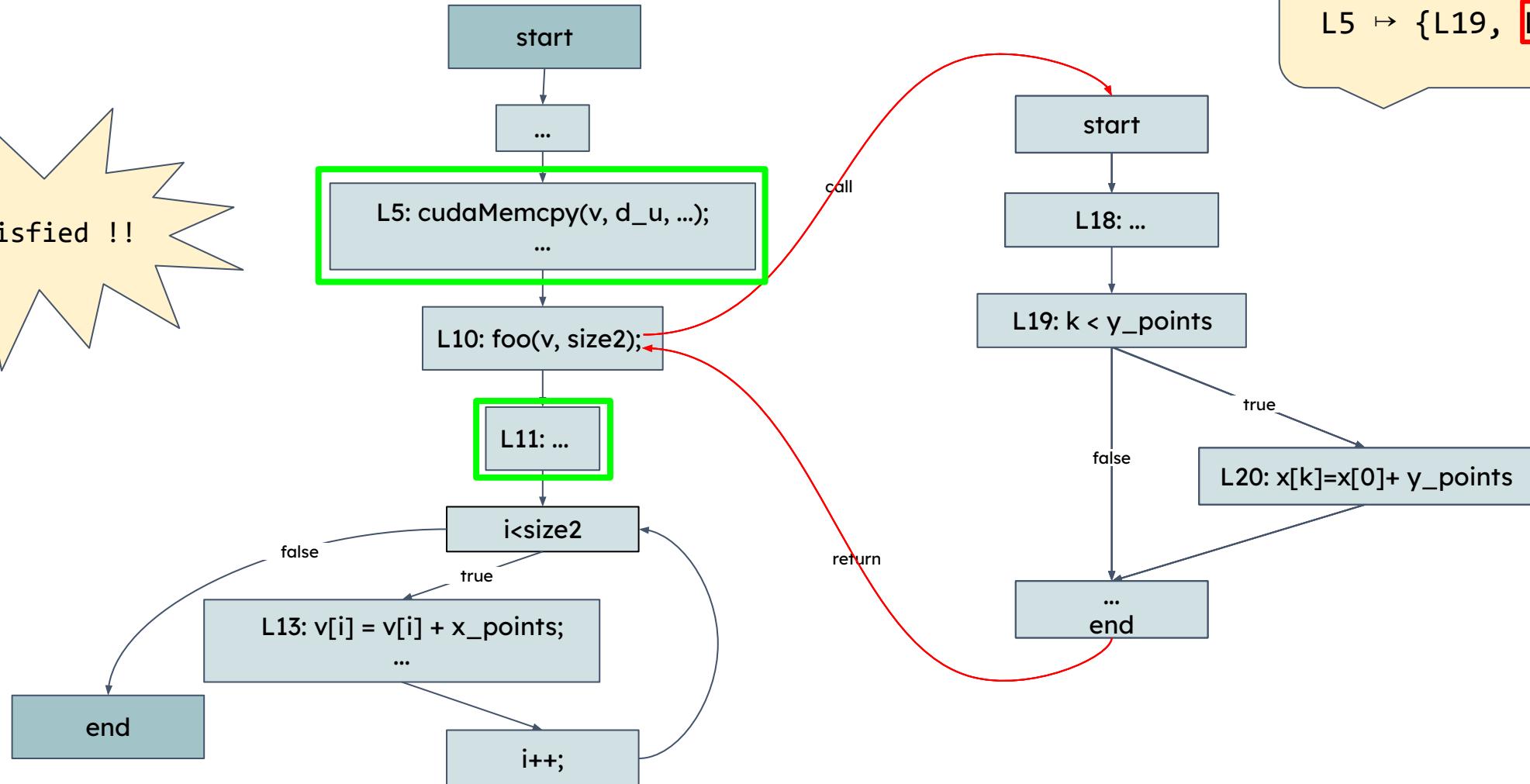


Step 2: Updating Target Locations (Contd.)

satisfied !!



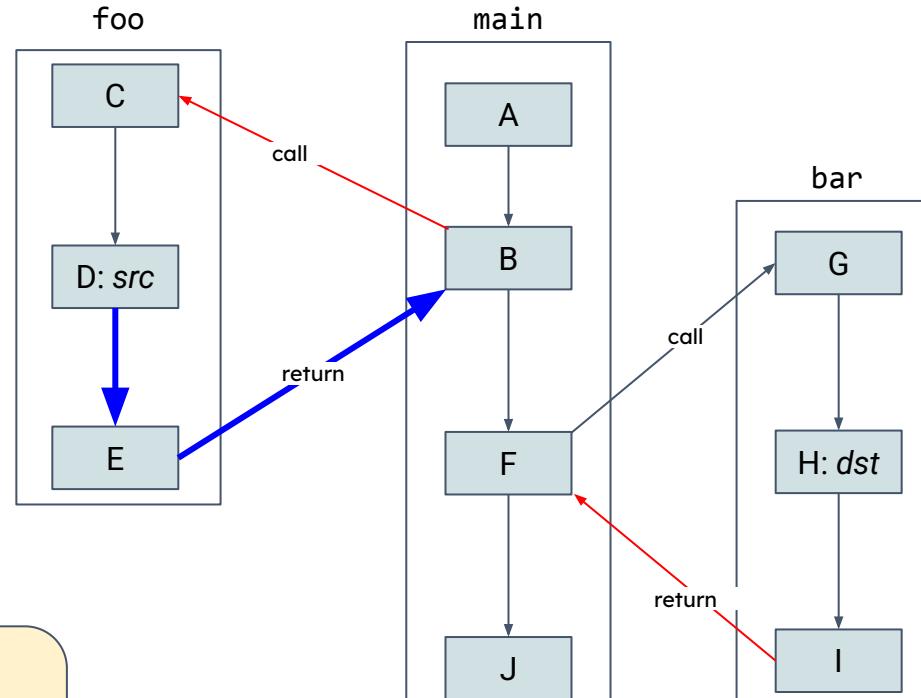
Step 2: Updating Target Locations (Contd.)



Target Map

$L5 \mapsto \{L19, L11\}$

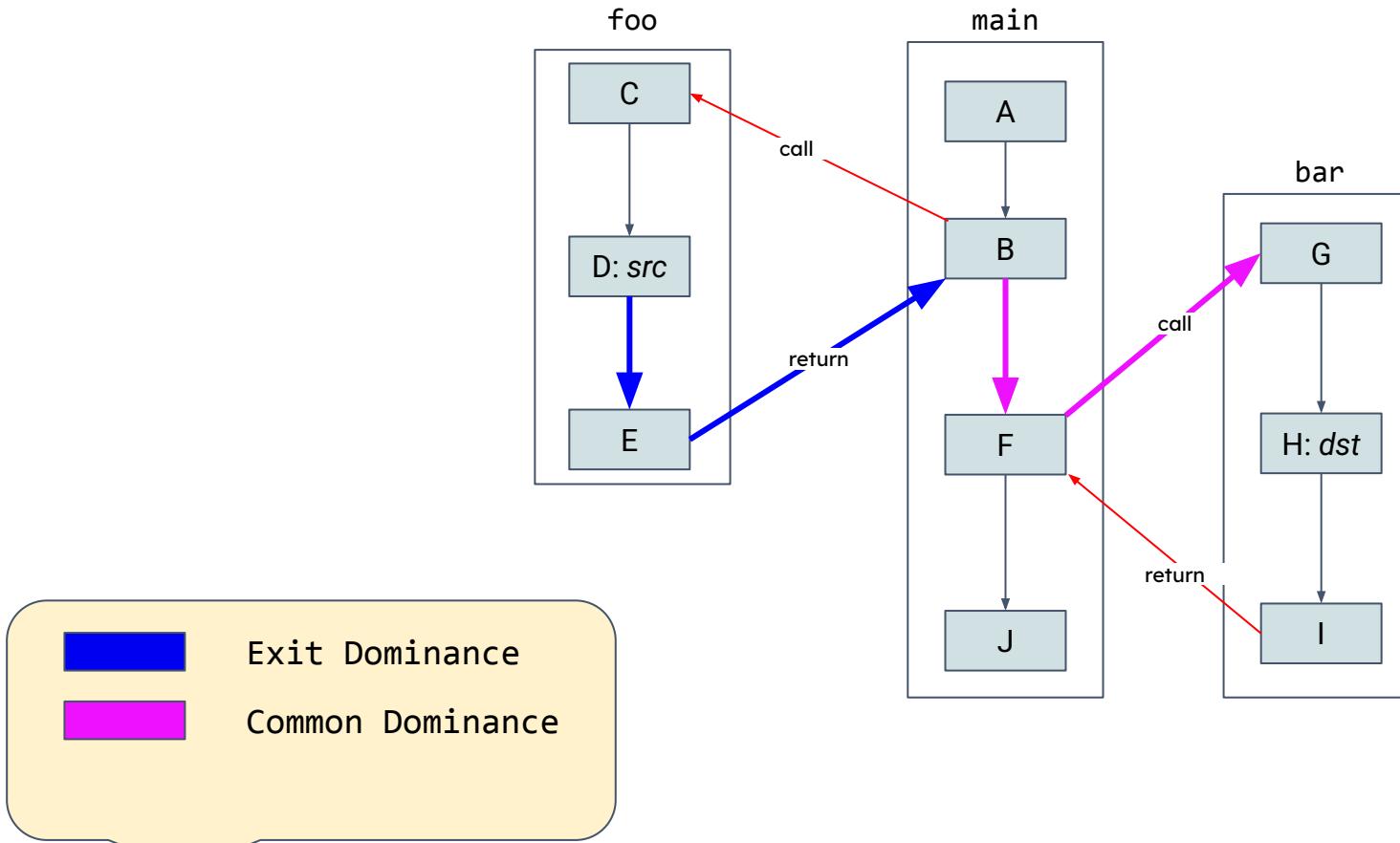
Step 2: Updating Target Locations: Inter-Procedural dom-pdom



- Demand-driven inter-procedural dom-pdom analysis
- Demand-driven query is $\langle D:\text{src}, H:\text{dst} \rangle$
- Based on transitivity of dom-pdom

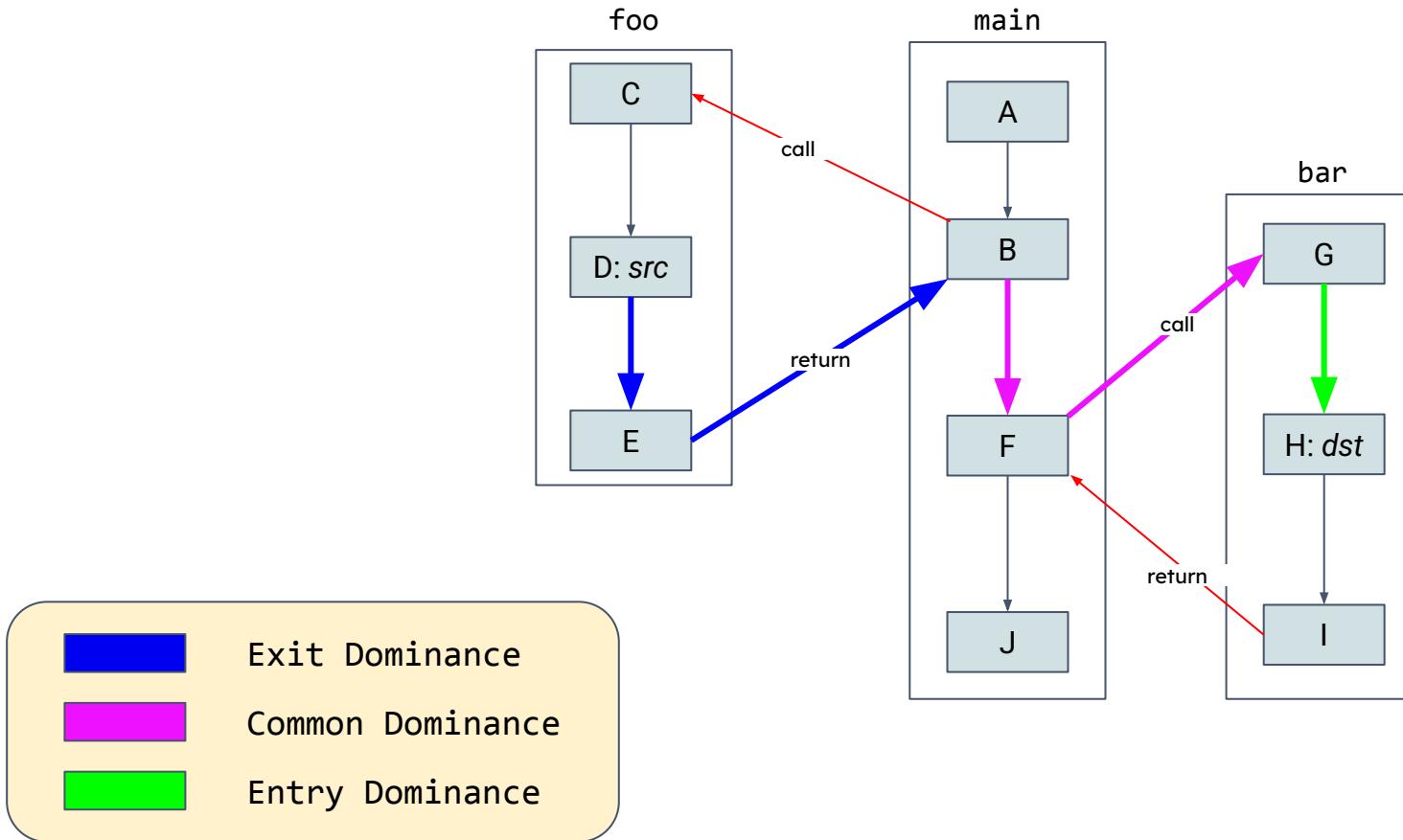
Exit Dominance

Step 2: Updating Target Locations: Inter-Procedural dom-pdom



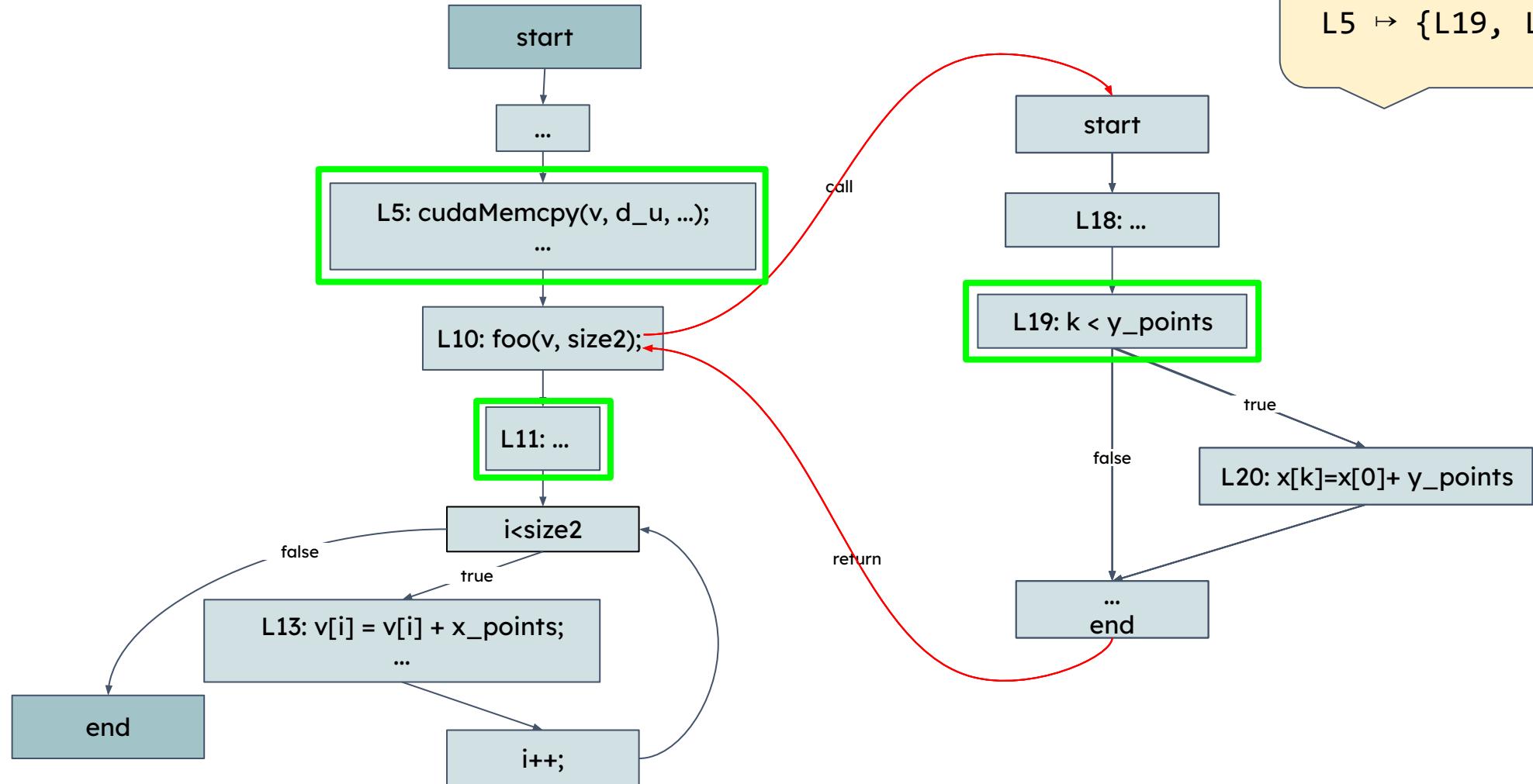
- Demand-driven inter-procedural dom-pdom analysis
- Demand-driven query is $\langle D:\text{src}, H:\text{dst} \rangle$
- Based on transitivity of dom-pdom

Step 2: Updating Target Locations: Inter-Procedural dom-pdom



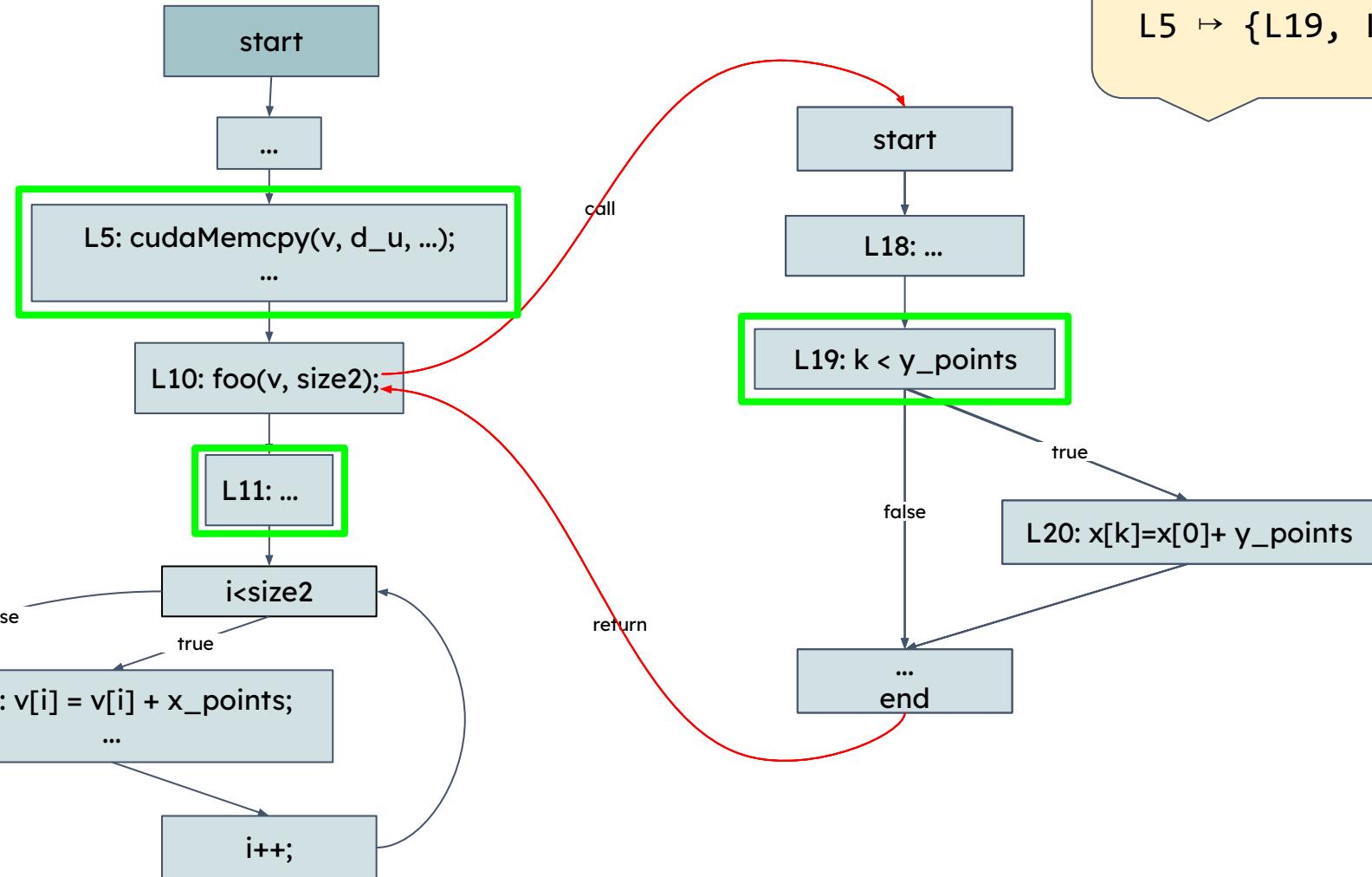
- Demand-driven inter-procedural dom-pdom analysis
- Demand-driven query is $\langle D:\text{src}, H:\text{dst} \rangle$
- Based on transitivity of dom-pdom

Step 3: Unifying Target Locations: Motivation



Step 3: Unifying Target Locations: Motivation

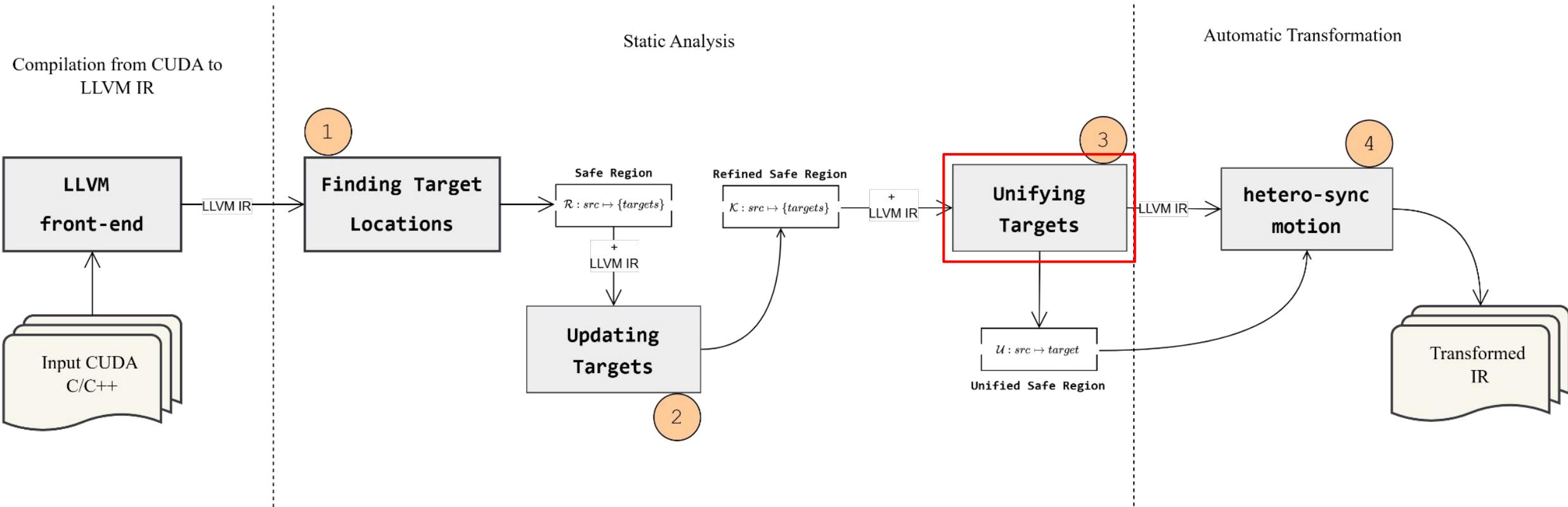
Memcpy statement
should not be
percolated to two
different locations



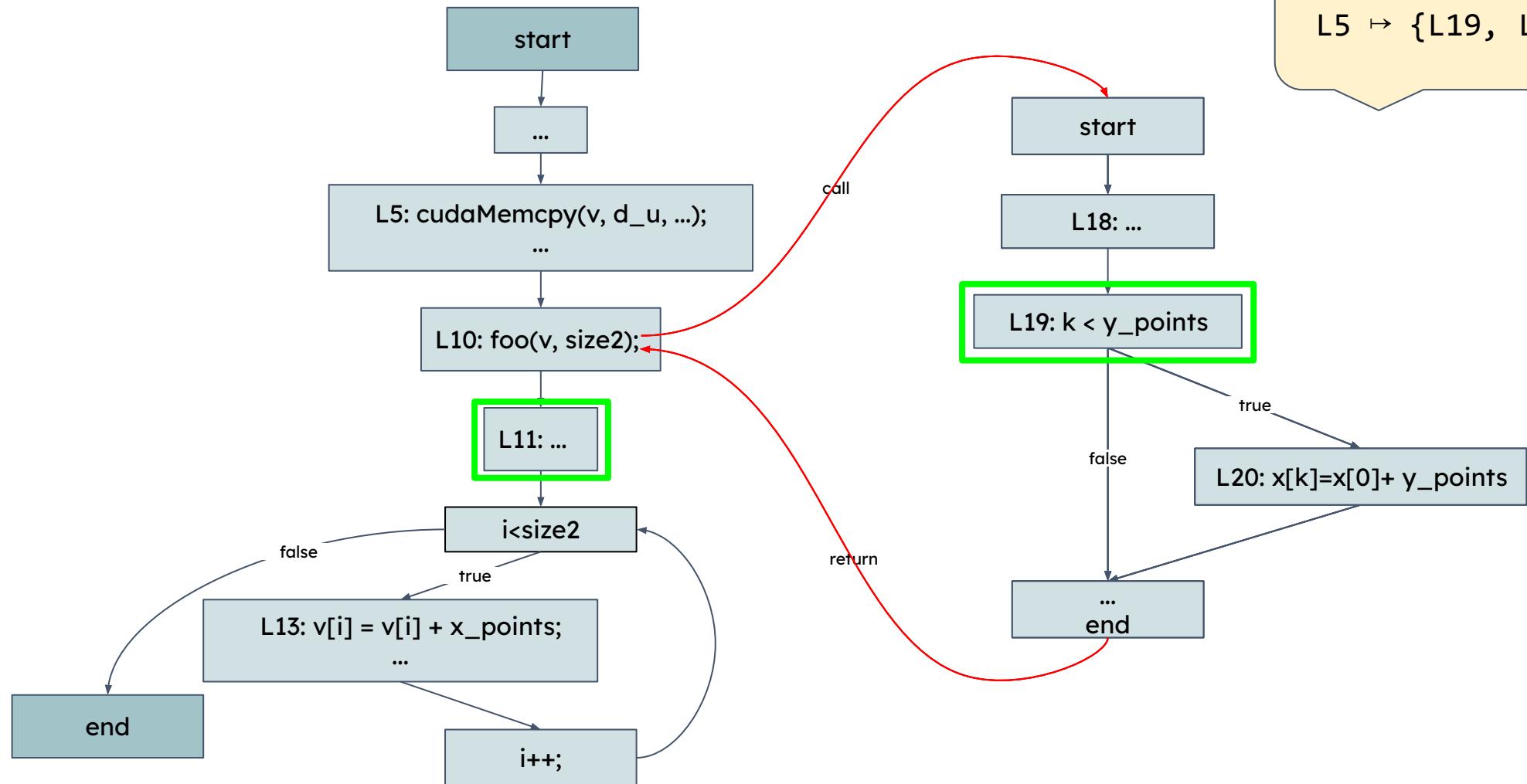
Target Map

$L5 \mapsto \{L19, L11\}$

Step 3: Unifying Target Locations

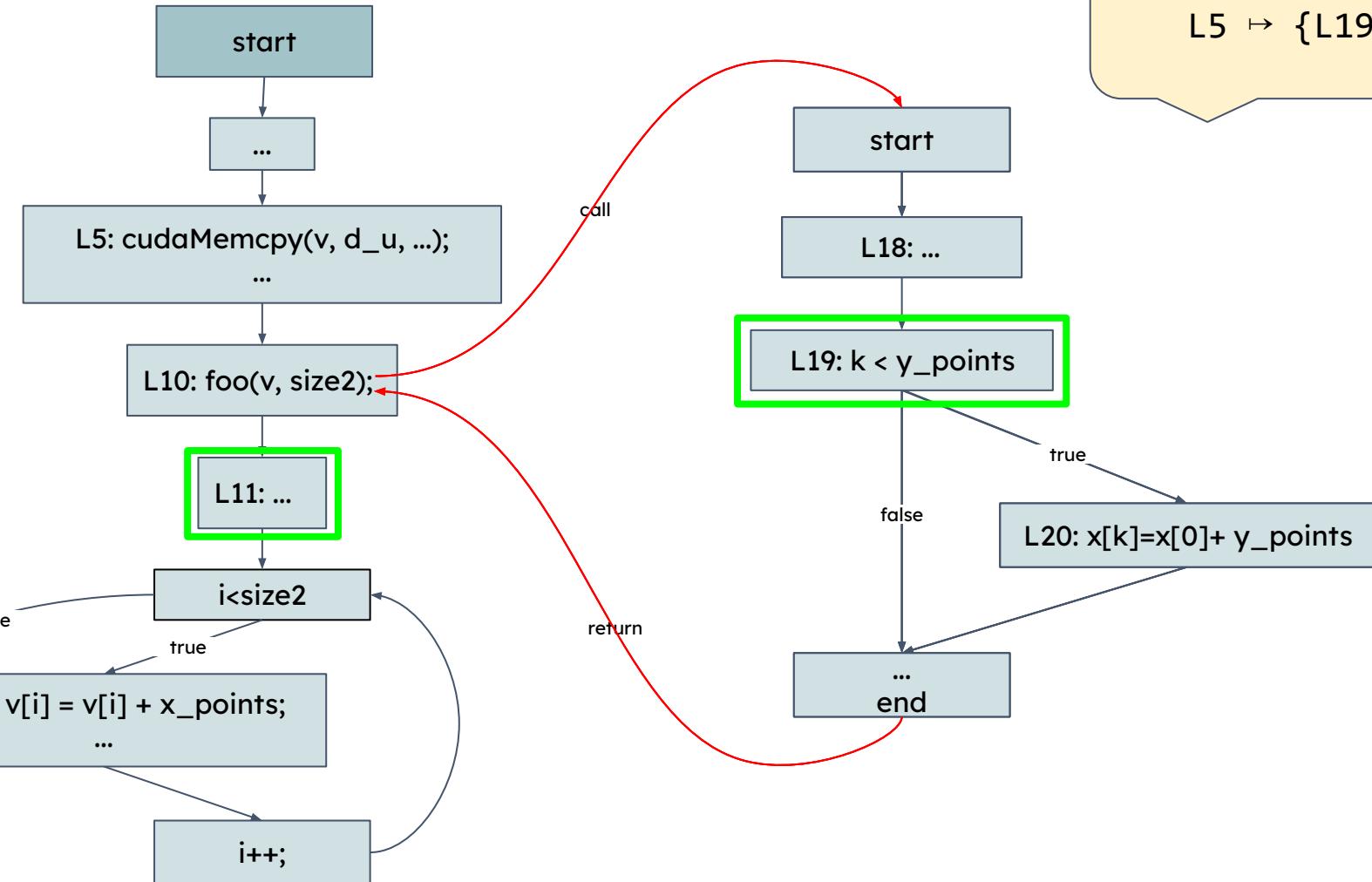


Step 3: Unifying Target Locations: Motivation



Step 3: Unifying Target Locations: Motivation

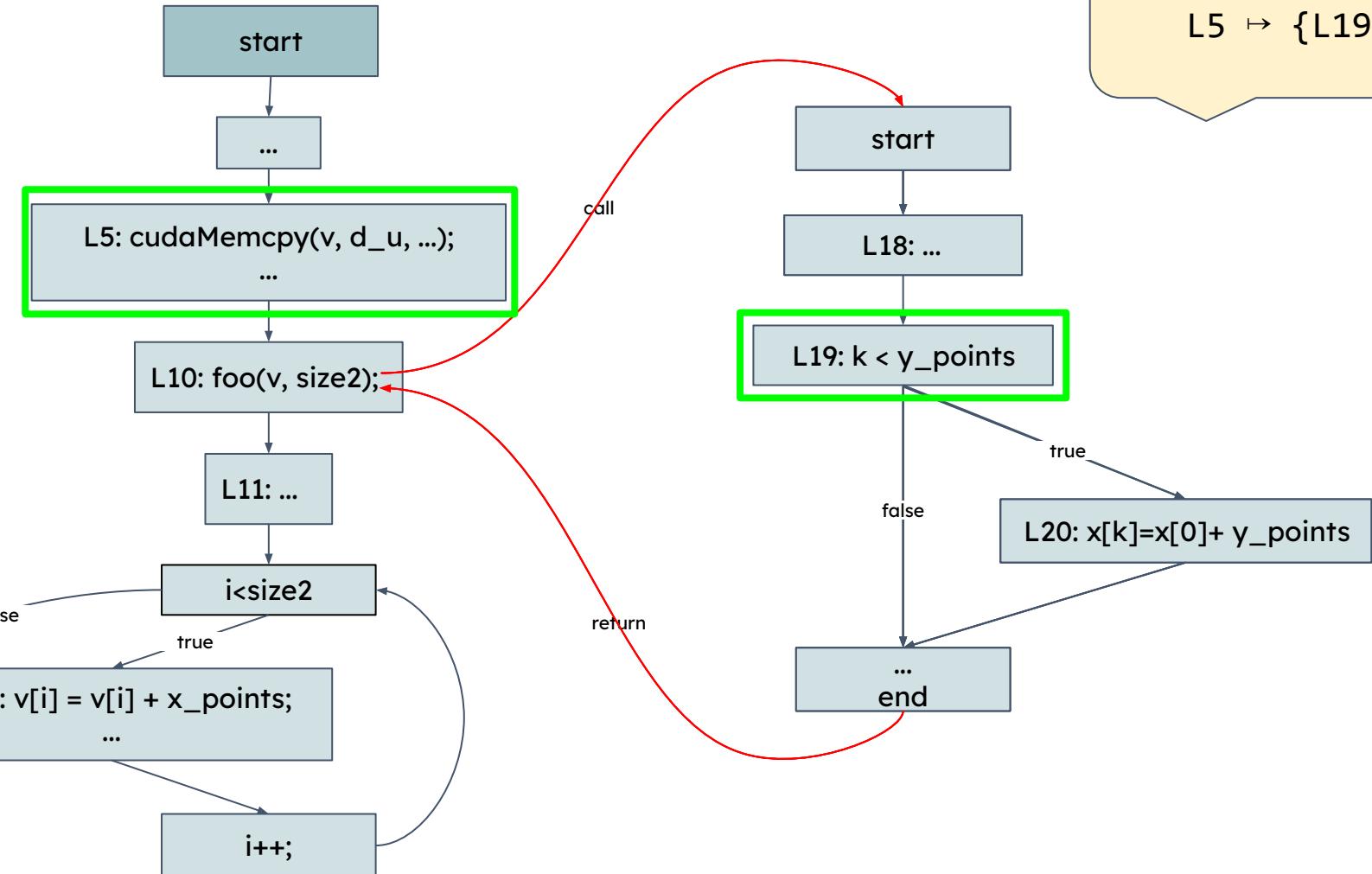
L19 dominates L11



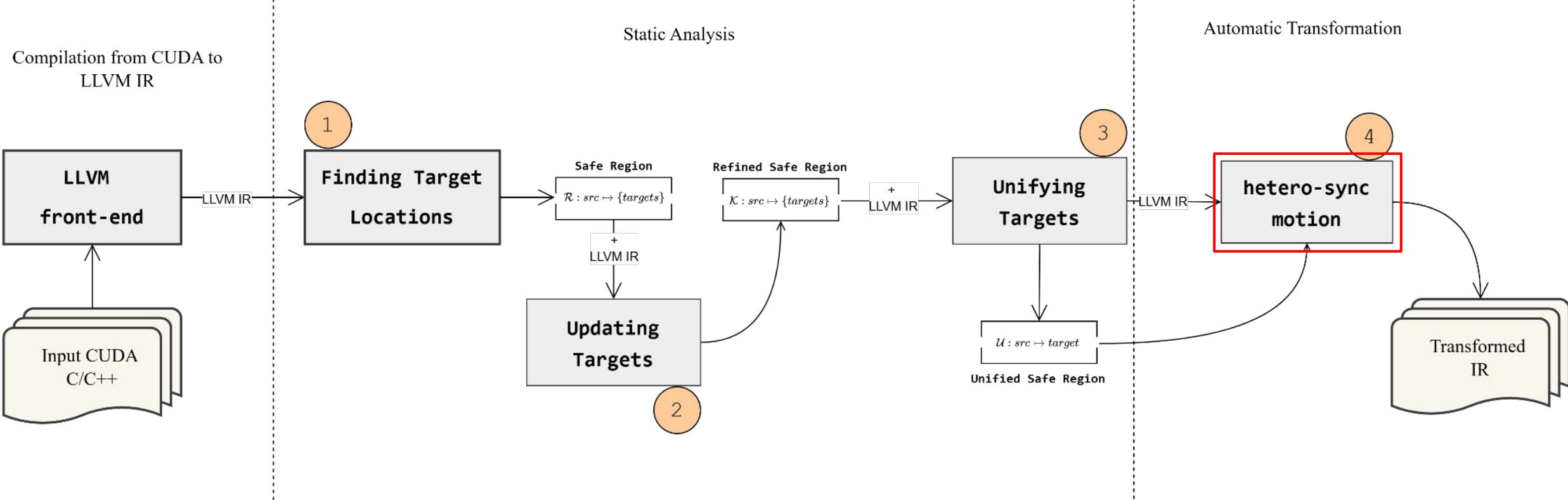
Target Map
 $L5 \mapsto \{L19\}$

Step 3: Unifying Target Locations: Motivation

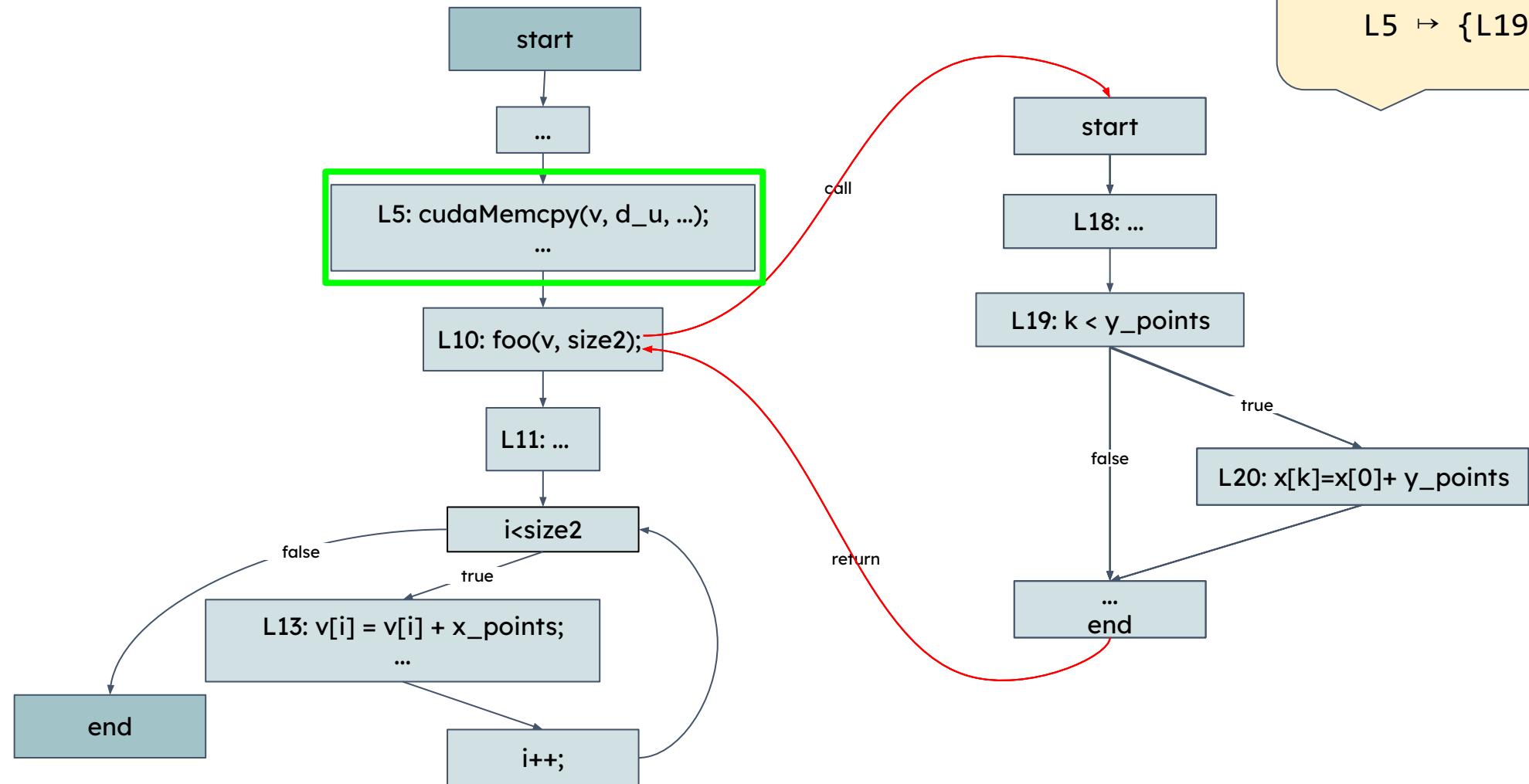
Finally we find one to one mapping between source and target



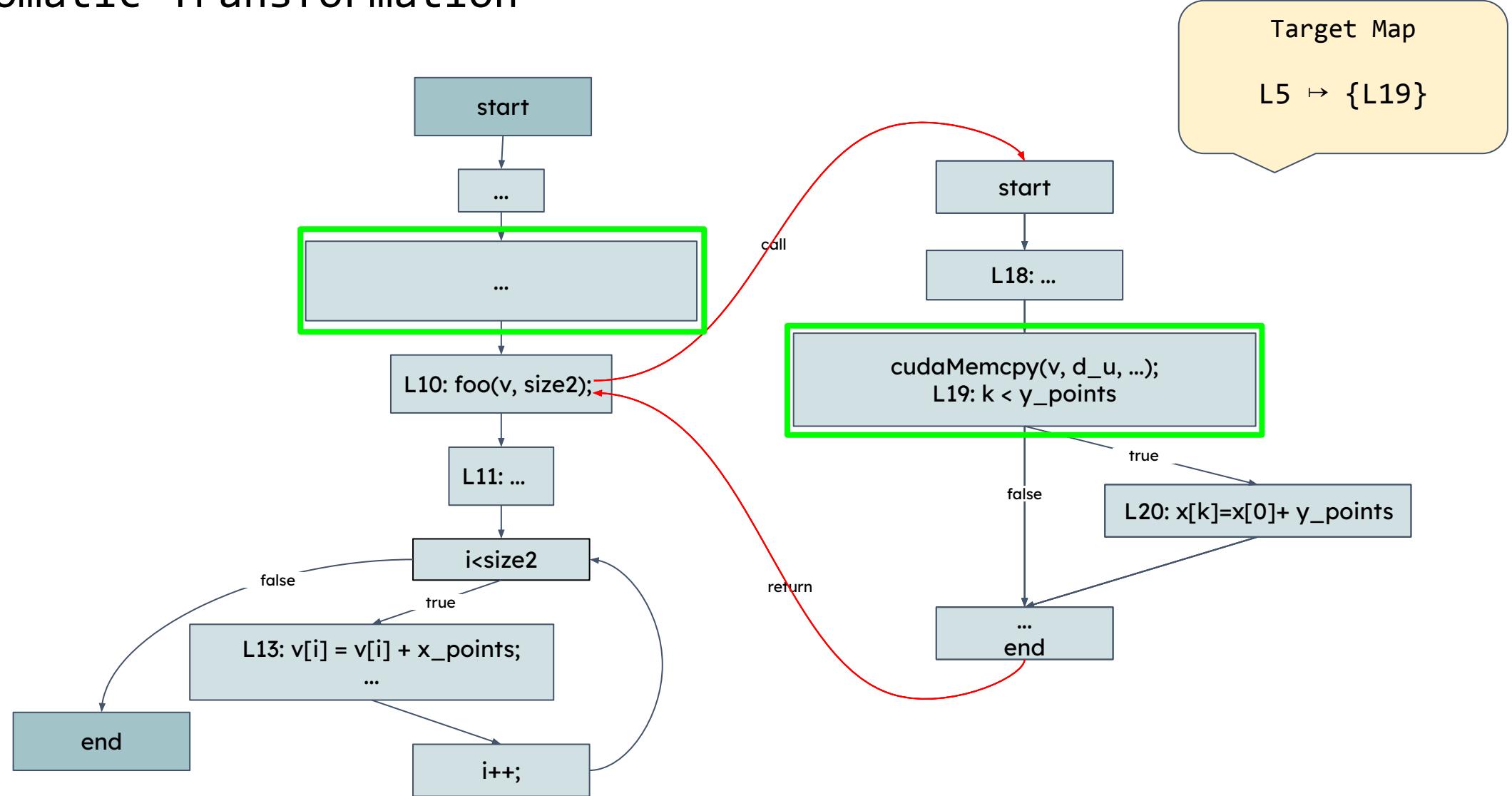
Automatic Transformation



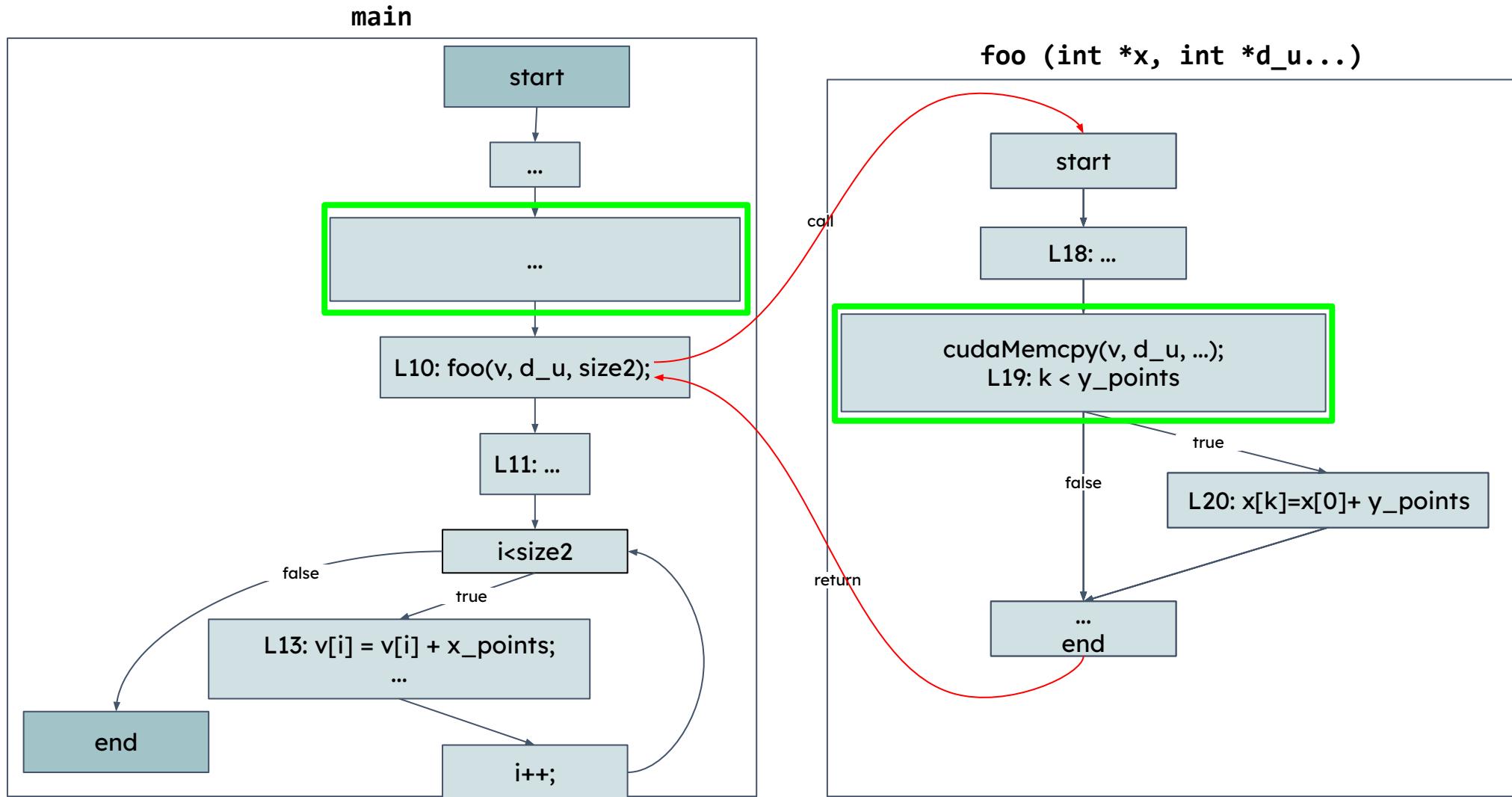
Automatic Transformation



Automatic Transformation



Automatic Transformation



Evaluation

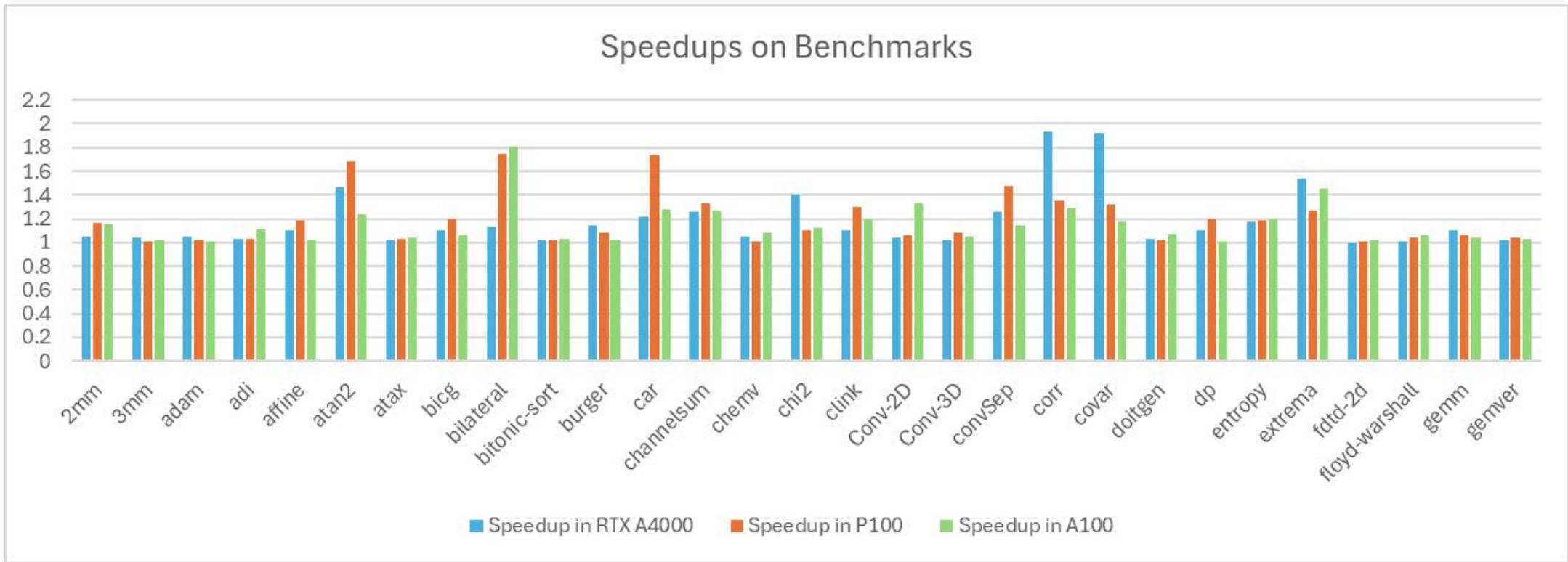
Evaluation: Experimental Setup

Component	Platform 1 (RTX A4000)	Platform 2 (P100)	Platform 3 (A100)
GPU (DRAM Size)	NVIDIA RTX A4000 (16 GB)	NVIDIA P100 (16 GB)	NVIDIA A100 (84 GB)
CPU	64 × AMD EPYC @ 2.20GHz	20 × Intel Xeon Silver @ 2.20GHz	64 × Intel Xeon Gold @ 3.50GHz
OS	Ubuntu 22.04 LTS	Debian GNU/Linux 11	Ubuntu 22.04 LTS
CPU RAM	128 GB	189 GB	84 GB

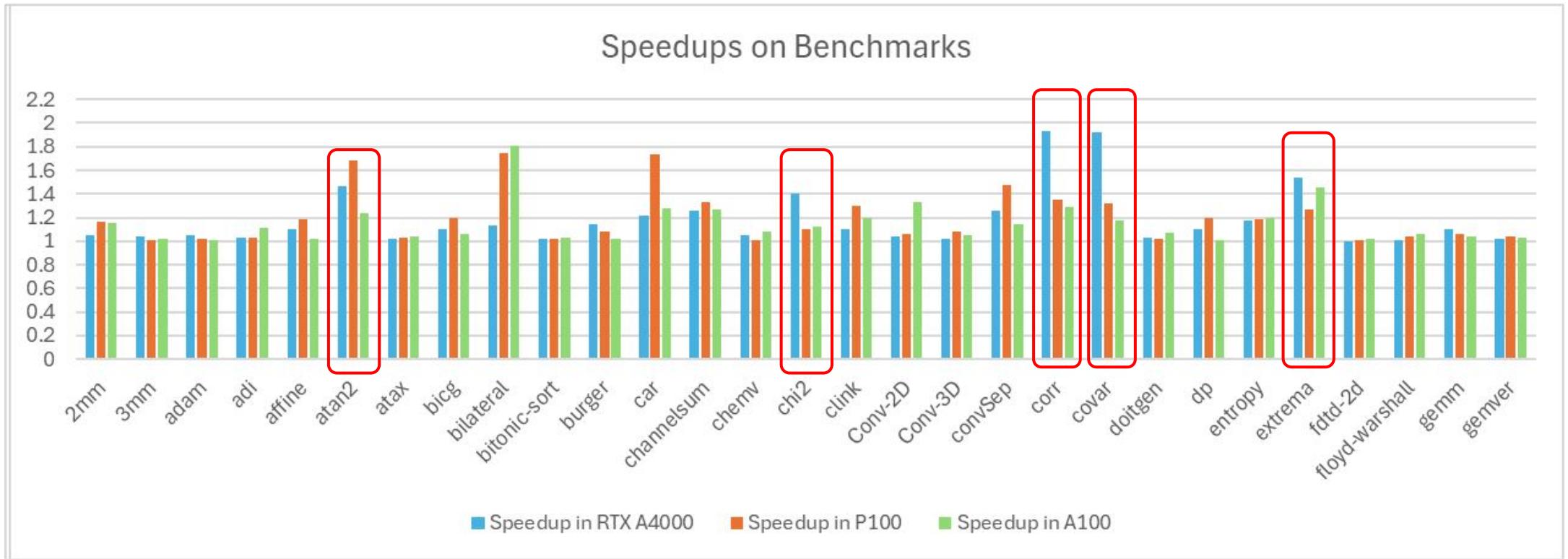
Evaluation: Efficiency of GSOHC

- Found **361** opportunities in **56** benchmark programs!!
- Around **48%** of the benchmarks have inter-procedural opportunities
- Benchmark LOC ranges from few hundreds to **18k**
- **Analysis time**
 - ranges from few milliseconds to 291 milliseconds
 - depends less on program size but more on control flow

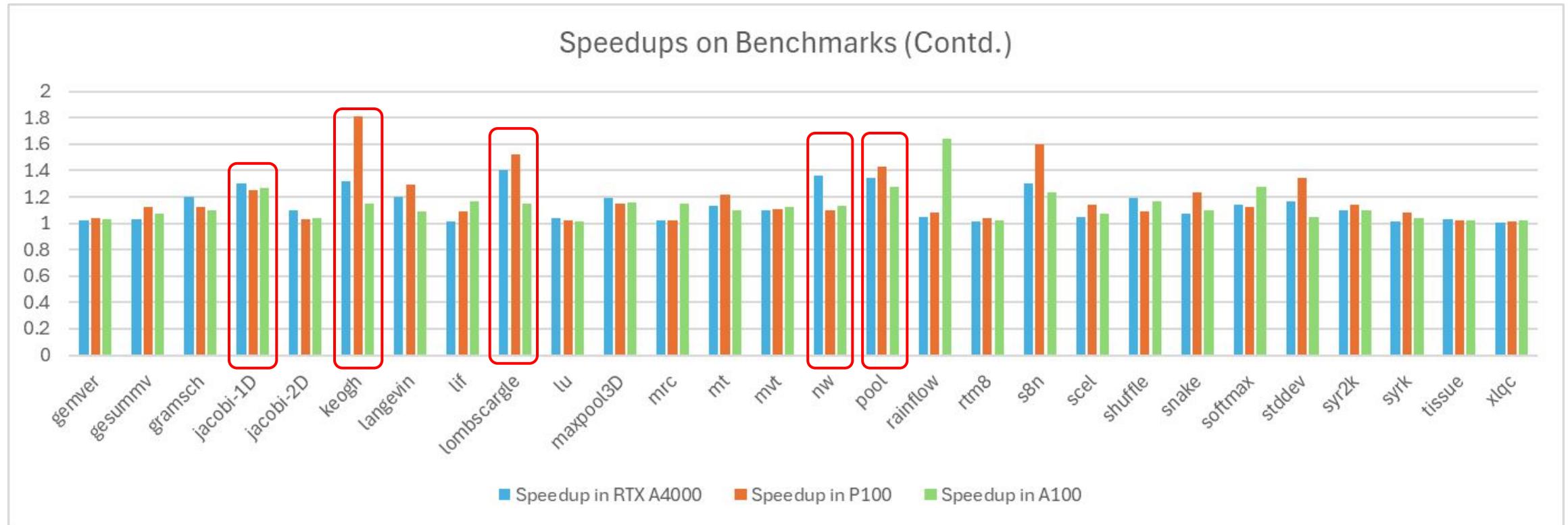
Evaluation: Effectiveness of GSOHC



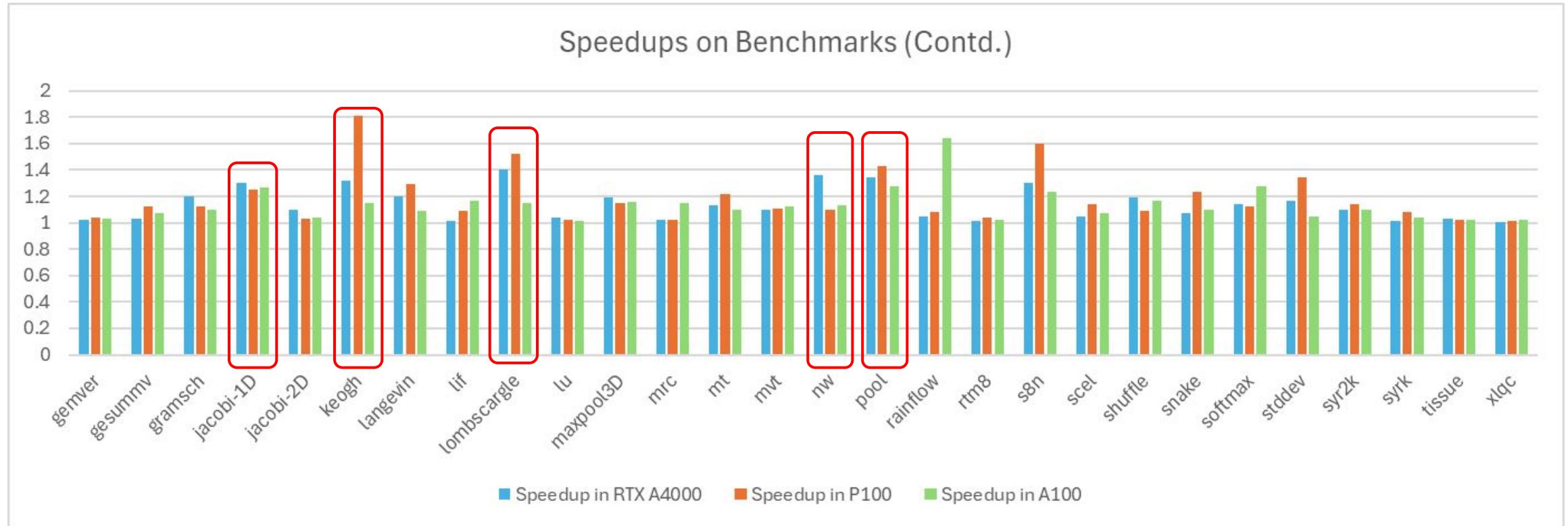
Evaluation: Effectiveness of GSOHC



Evaluation: Effectiveness of GSOHC (Contd.)

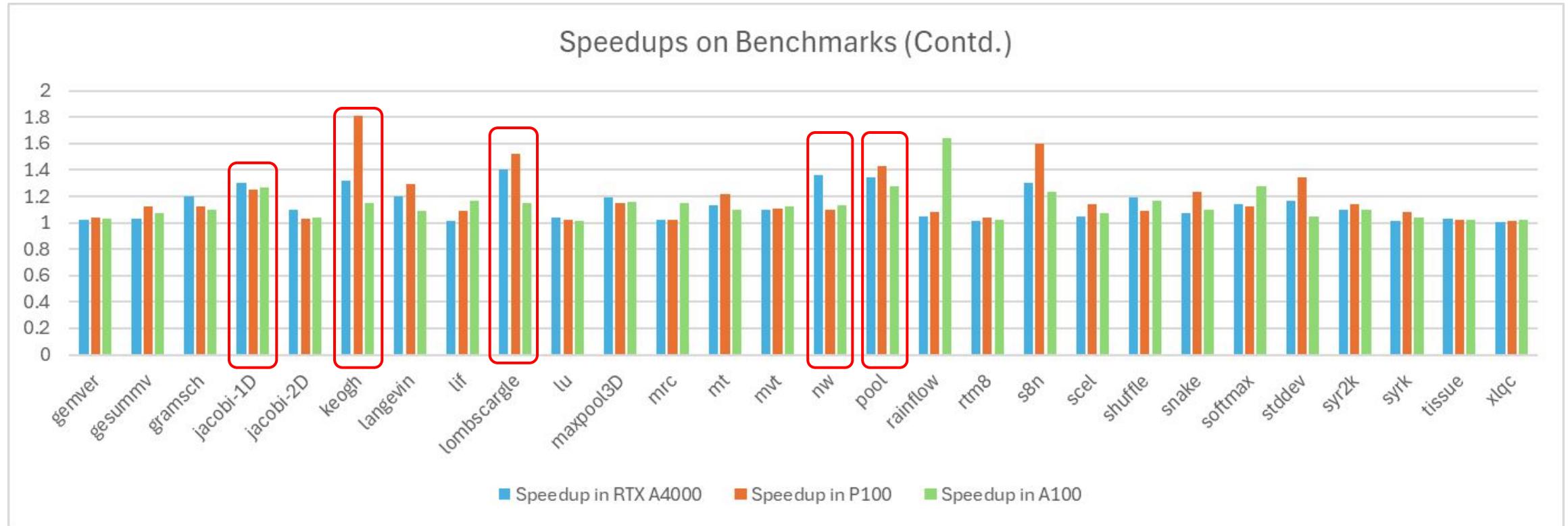


Evaluation: Effectiveness of GSOHC



Performance Improvement: Speedup upto 1.9x and 34% improvement with an average of 11%

Evaluation: Effectiveness of GSOHC



Insight: The optimization yields best performance if the host-side computation time, and the device-side computation time, are significant and are comparable to one another.

Summary

- **Hetero-Sync Motion:** Moves synchronization statements to improve CPU-GPU overlap → better performance.
- **GSOHC**
 - A novel compiler framework for optimizing CPU-GPU synchronization
 - 3-phased framework
 - Interprocedural, context- and flow-sensitive data-flow analysis
- **Comprehensive Evaluation**
 - Tested on 56 benchmark programs
 - Identified all optimization opportunities
 - Achieved up to 1.9× speedup and 34% improvement with an average of 11%
 -

Summary

- **Hetero-Sync Motion:** Moves synchronization statements to improve CPU-GPU overlap → better performance.
- **GSOHC**
 - A novel compiler framework for optimizing CPU-GPU synchronization
 - 3-phased framework
 - Interprocedural, context- and flow-sensitive data-flow analysis
- **Comprehensive Evaluation**
 - Tested on 56 benchmark programs
 - Identified all optimization opportunities
 - Achieved up to 1.9× speedup

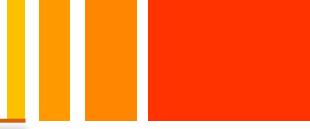
Webpage



Artifact



Thank You!!



Webpage



LinkedIn

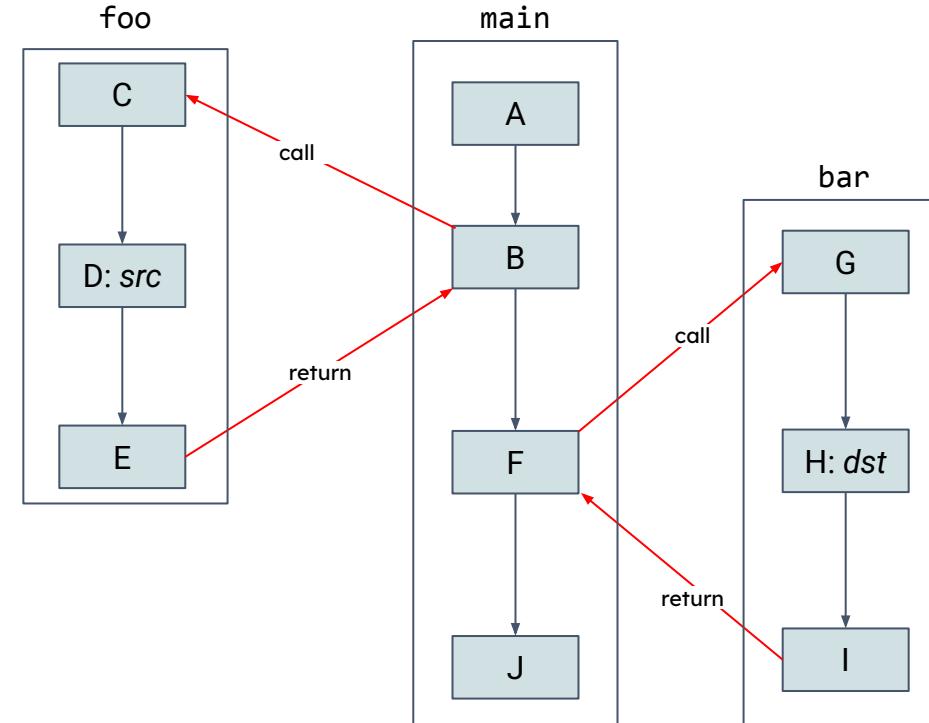


Artifact

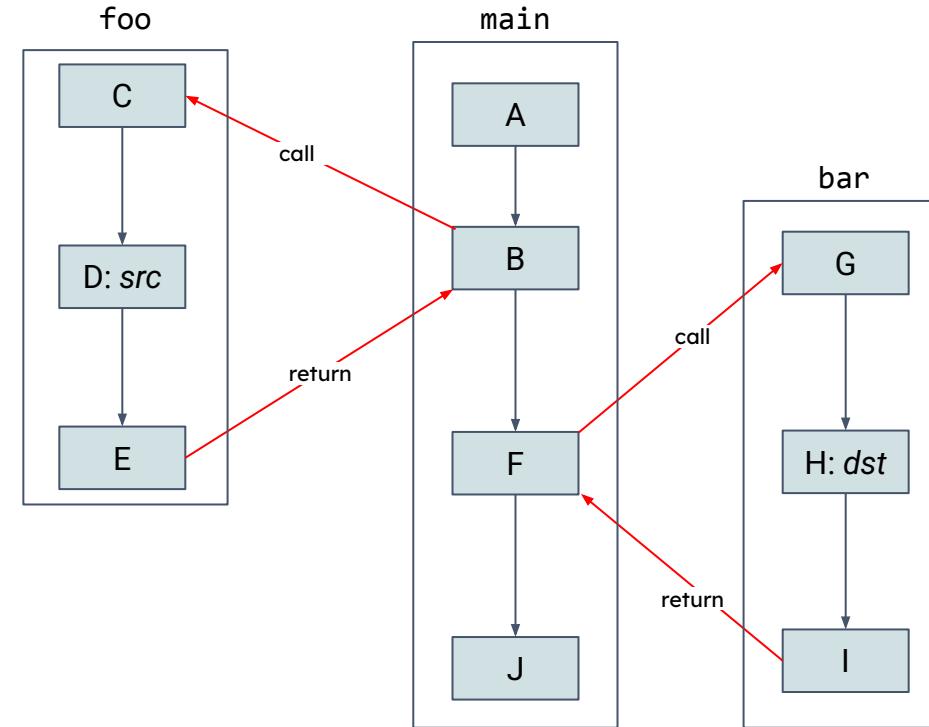


Thank You

Step 2: Updating Target Locations: Inter-Procedural dom - pdom

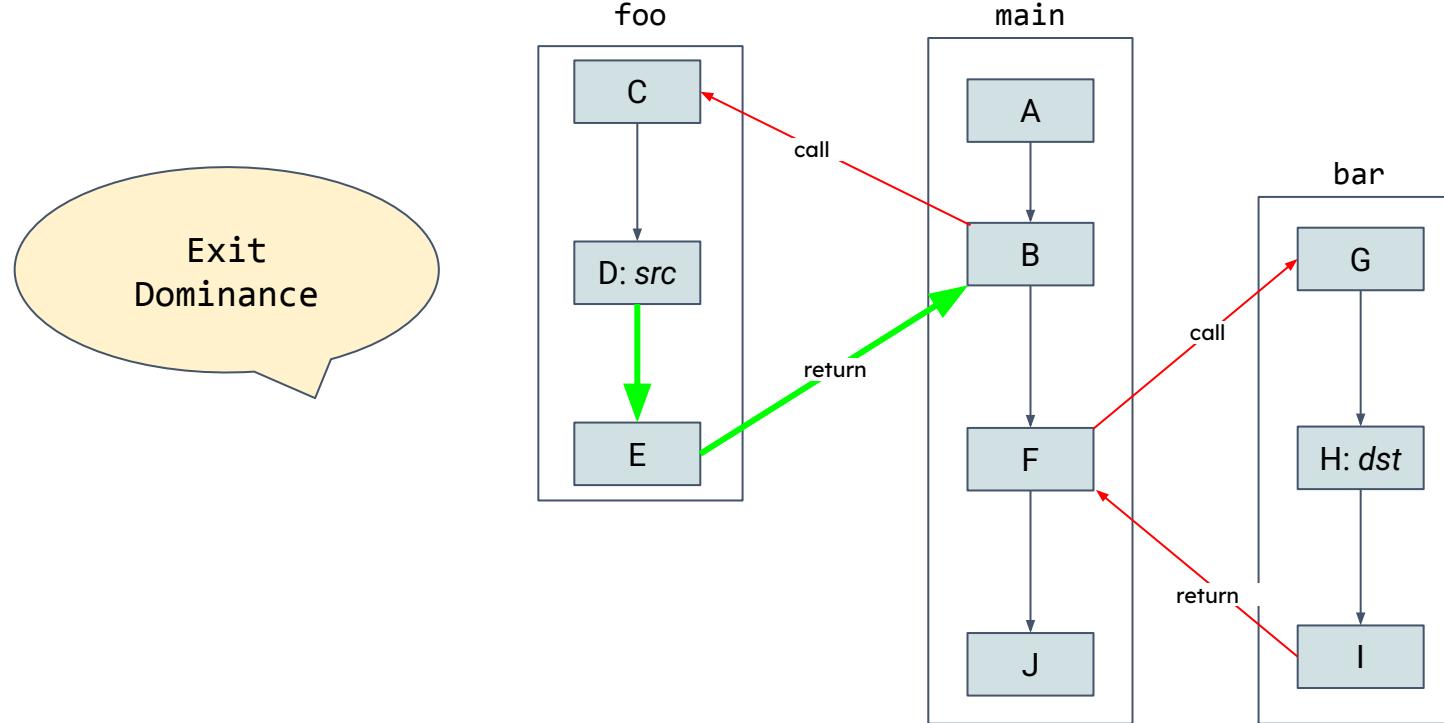


Step 2: Updating Target Locations: Inter-Procedural dom - pdom

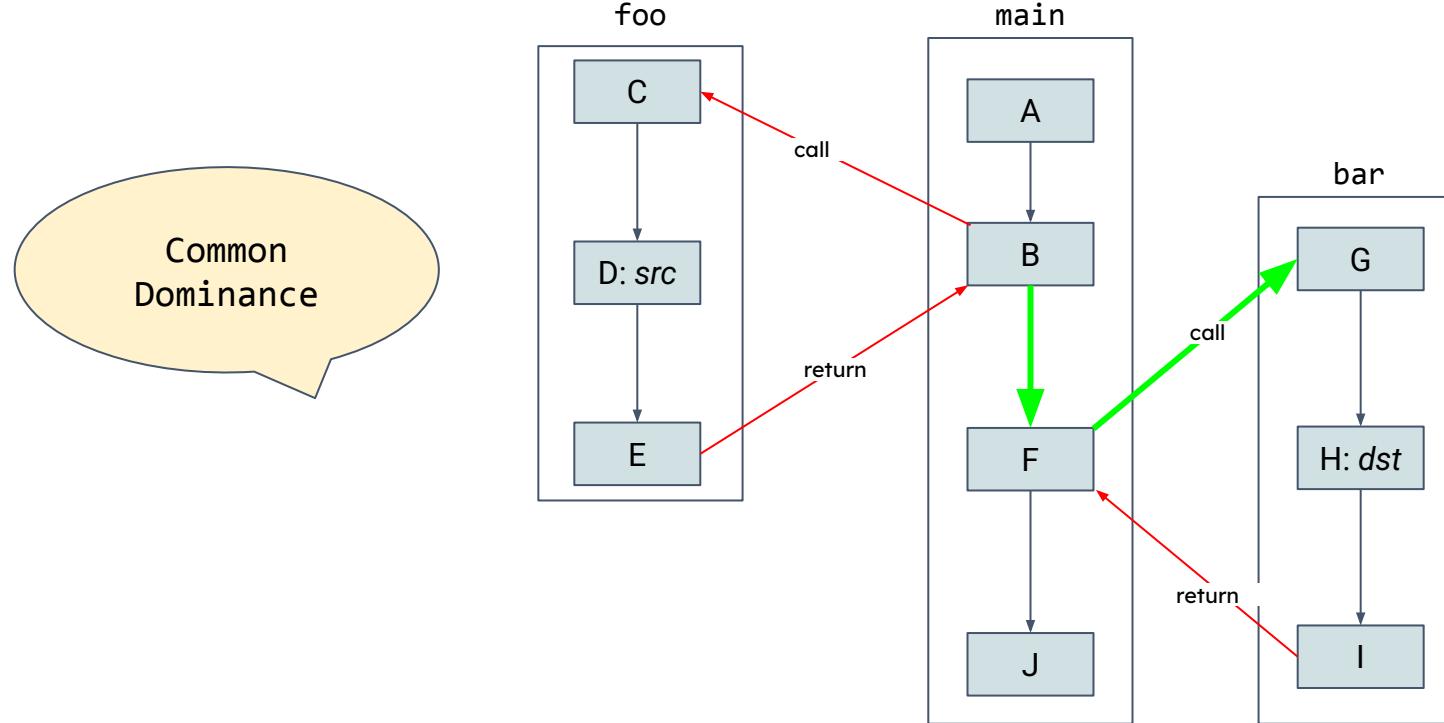


src *dom* dst and
dst *pdom* src

Step 2: Updating Target Locations: Inter-Procedural dom - pdom



Step 2: Updating Target Locations: Inter-Procedural dom - pdom



Step 2: Updating Target Locations: Inter-Procedural dom - pdom

