

CP for Bin Packing with Multi-core and GPUs

Fabio Tardivo¹ Laurent Michel² Enrico Pontelli¹

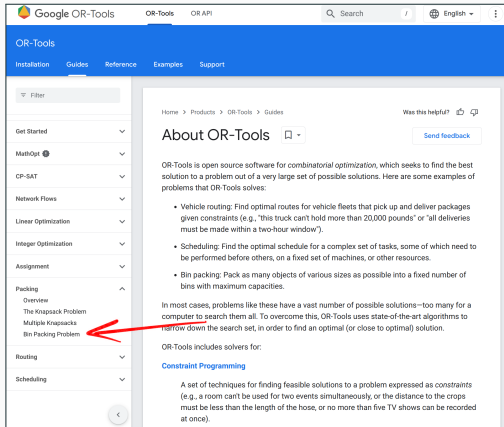
¹New Mexico State University

²University of Connecticut



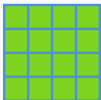
Motivation

Packing items is a crucial problem in many applications.

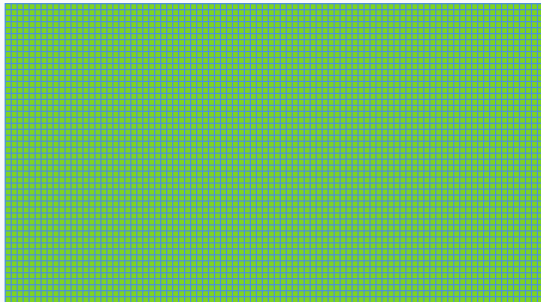


The screenshot shows the Google OR-Tools website. The navigation menu on the left includes: Installation, Guides, Reference, Examples, and Support. Under the 'Guides' section, there is a list of topics: Get Started, MathOpt, CP-SAT, Network Flows, Linear Optimization, Integer Optimization, Assignment, Packing, Routing, and Scheduling. The 'Packing' section is expanded, showing sub-items: Overview, The Knapsack Problem, Multiple Knapsacks, and Bin Packing Problem. A red arrow points to the 'Bin Packing Problem' link. The main content area displays the 'About OR-Tools' page, which includes a list of examples of problems solved by OR-Tools: Vehicle routing, Scheduling, and Bin packing. The Bin packing description states: 'Pack as many objects of various sizes as possible into a fixed number of bins with maximum capacities.'

Cores difference between CPU and GPU.

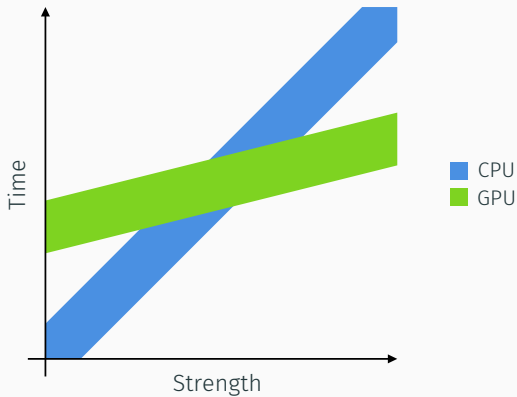


CPU

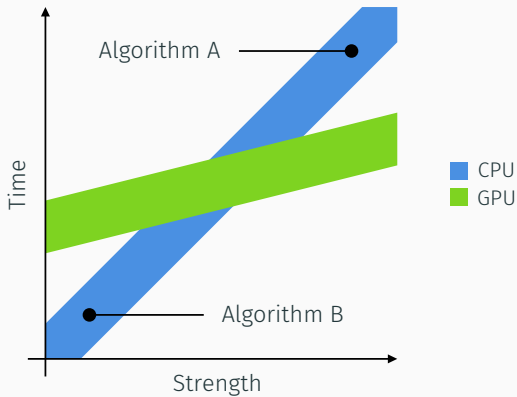


GPU

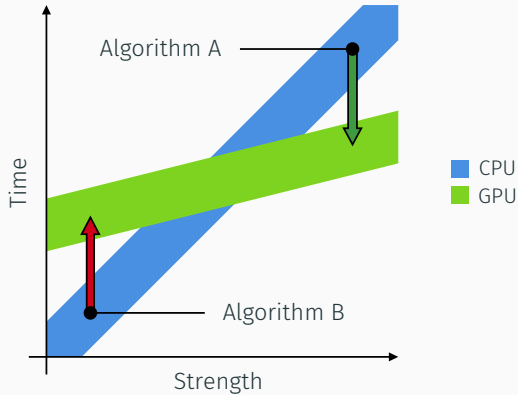
Different trade-off propagation strength vs propagation time.



Different trade-off propagation strength vs propagation time.

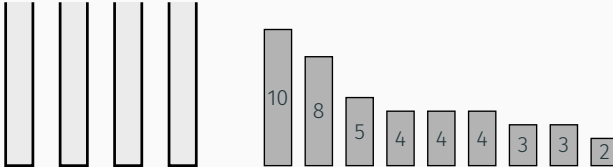


Different trade-off propagation strength vs propagation time.

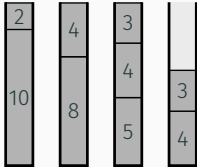


Introduction

The Bin Packing Problem (BPP) involves packing n items with weights $W = [w_1, w_2, \dots, w_n]$ into the fewest bins of capacity c .



The Bin Packing Problem (BPP) involves packing n items with weights $W = [w_1, w_2, \dots, w_n]$ into the fewest bins of capacity c .



An instance with n items of weights $W = [w_1, w_2, \dots, w_n]$, and k bins of capacity c is modeled as:

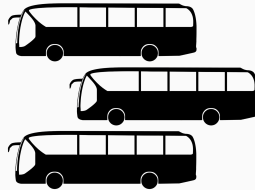
$$x_i = \{1, \dots, k\} \quad i = 1, \dots, n$$

$$l_j = \{0, \dots, c\} \quad j = 1, \dots, k$$

$$\text{BinPacking}(W = [w_1, \dots, w_n], X = [x_1, \dots, x_n], L = [l_1, \dots, l_k])$$

Simplified propagation algorithm for the BinPacking constraint:

```
Procedure: propagate( $c, W, k, X, L$ )  
for  $j \leftarrow 1$  to  $k$  do  
  doLoadCoherence( $j, X, W, L$ )  
  doLoadTightening( $j, X, W, L$ )  
  for  $i \in \{i \mid j \in x_i \wedge |x_i| > 1\}$  do  
    doItemEliminationCommitment( $i, j, X, W, L$ )  
if getLowerBound( $c, W, k, X$ )  $> k$  then Fail
```



- Fast.
- Few people/tasks.

- Slow.
- Many people/tasks.

GPU execution model is known as Single-Instruction Multiple-Thread:

- Each thread executes the same C/C++ function, called a kernel.
- Each thread has a unique index used for data access and control flow.

GPU execution model is known as Single-Instruction Multiple-Thread:

- Each thread executes the same C/C++ function, called a kernel.
- Each thread has a unique index used for data access and control flow.

Sequential combination of the arrays A and B into C :

Procedure *Main*(Args)

```
n ← ...  
A ← ...  
... // Initialization  
  
combineArrays(n, A, B, C)
```

Procedure *combineArrays*(n, A, B, C)

```
for  $i \leftarrow 0$  to  $n - 1$  do  
  if  $i \bmod 2 = 0$  then  
    |  $C[i] \leftarrow C[i] + f(A[i], B[i])$   
  else  
    |  $C[i] \leftarrow C[i] + g(A[i], B[i])$ 
```

GPU execution model is known as Single-Instruction Multiple-Thread:

- Each thread executes the same C/C++ function, called a kernel.
- Each thread has a unique index used for data access and control flow.

Parallel combination of the arrays *A* and *B* into *C*:

Procedure *Main*(*Args*)

```
n ← ...
A ← ...
... // Initialization

nThreads ← n
memcpyCpuToGpuAsync(n, A, B, C)
launchKernelAsync(combineArraysKernel, nThreads, [n, A, B, C])
memcpyGpuToCpuAsync(C)
waitGpu()
```

Procedure *combineArraysKernel*(*n*, *A*, *B*, *C*)

```
i ← getThreadId() // n threads
if i mod 2 = 0 then
  | C[i] ← C[i] + f(A[i], B[i])
else
  | C[i] ← C[i] + g(A[i], B[i])
```


GPU execution model is known as Single-Instruction Multiple-Thread:

- Each thread executes the same C/C++ function, called a kernel.
- Each thread has a unique index used for data access and control flow.

Parallel combination of the arrays A and B into C:

}}}}}}}}}}}}}}}}}}}}}}}}

```
Procedure combineArraysKernel(n, A, B, C)
  i ← getThreadIdx() // n threads
  if i mod 2 = 0 then
    C[i] ← C[i] + f(A[i], B[i])
  else
    C[i] ← C[i] + g(A[i], B[i])
```

GPU execution model is known as Single-Instruction Multiple-Thread:

- Each thread executes the same C/C++ function, called a kernel.
- Each thread has a unique index used for data access and control flow.

Parallel combination of the arrays A and B into C :



```
Procedure combineArraysKernel( $n, A, B, C$ )  
   $i \leftarrow \text{getThreadIdx}() // n \text{ threads}$   
  if  $i \bmod 2 = 0$  then  
     $C[i] \leftarrow C[i] + f(A[i], B[i])$   
  else  
     $C[i] \leftarrow C[i] + g(A[i], B[i])$ 
```

GPU execution model is known as Single-Instruction Multiple-Thread:

- Each thread executes the same C/C++ function, called a kernel.
- Each thread has a unique index used for data access and control flow.

Parallel combination of the arrays *A* and *B* into *C*:

Procedure *Main*(*Args*)

```
n ← ...  
A ← ...  
... // Initialization
```

```
nThreads ←  $\frac{n}{2}$   
memcpyCpuToGpuAsync(n, A, B)  
launchKernelAsync(combineArraysKernelEven, nThreads, ...)  
launchKernelAsync(combineArraysKernelOdd, nThreads, ...)  
memcpyGpuToCpuAsync(C)  
waitGpu()
```

Procedure *combineArraysKernelEven*(*n*, *A*, *B*, *C*)

```
i ← getThreadId() // n/2 threads  
j ← 2i  
C[j] ← C[i] + f(A[i], B[i])
```

Procedure *combineArraysKernelOdd*(*n*, *A*, *B*, *C*)

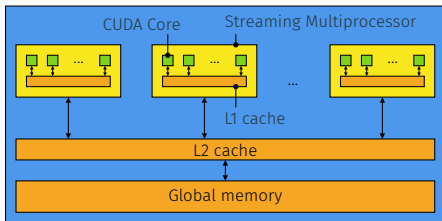
```
i ← getThreadId() // n/2 threads  
j ← 1 + 2i  
C[j] ← C[i] + g(A[i], B[i])
```

Simplified GPU architecture:

- Streaming Multiprocessors
- CUDA Cores
- L1 Cache
- L2 Cache
- Global Memory

Memory latency in clock cycles:

- L1 cache: 30
- L2 cache: 260 (8x slower)
- Global memory: 470 (15x slower)



Atomic operations:

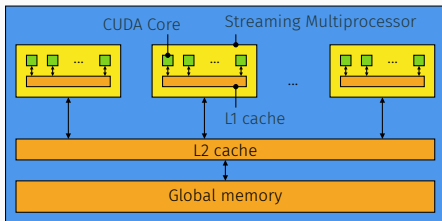
- Min, Max, And, Or, ...
- Block other access to the memory

NVIDIA GeForce RTX 3080 specifications:

- Streaming Multiprocessors: 68
 - CUDA Cores: 128 at 1.4 GHz
 - L1 Cache: 128 KB
 - L2 Cache: 5 MB
 - Global Memory: 10 GB
- } Total of **8704** threads!

Memory latency in clock cycles:

- L1 cache: 30
- L2 cache: 260 (8x slower)
- Global memory: 470 (15x slower)



Atomic operations:

- Min, Max, And, Or, ...
- Block other access to the memory

GPU acceleration enhances a wide range of fields:

- Machine Learning: PyTorch, TensorFlow, cuDNN, ...
- Numerical Analysis: MATLAB, cuBLAS, cuFFT, ...
- Scientific Simulation: GROMACS, ANSYS Fluent, Qiskit, ...
- Constraint Programming: **None**.

Constraint Programming has specific requirements, so it is necessary:

- Build **tailored** software from scratch.
- Make it **transparent** to the user.

Design

Simplified propagation algorithm for the BinPacking constraint

```
Procedure: propagate( $c, W, k, X, L$ )  
for  $j \leftarrow 1$  to  $k$  do  
  doLoadCoherence( $j, X, W, L$ )  
  doLoadTightening( $j, X, W, L$ )  
  for  $i \in \{i \mid j \in x_i \wedge |x_i| > 1\}$  do  
    doItemEliminationCommitment( $i, j, X, W, L$ )  
if getLowerBound( $c, W, k, X$ )  $> k$  then Fail
```


Two prominent approaches to obtain a tighter lower bound:

- Linear relaxation of the strong Arc-Flow model.
- Enhance the weak L_1 bound with a Dual Feasible Function (DFF).

Given an instance with weights W and capacity c , the L_1 bound is:

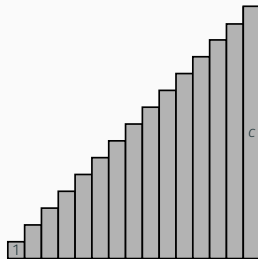
$$L_1(c, W) = \left\lceil \frac{1}{c} \sum_{w \in W} w \right\rceil$$

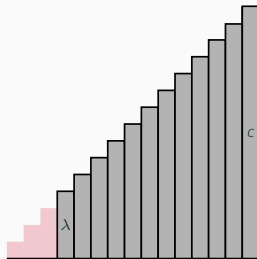
A function $f: \mathbb{N}_0 \rightarrow \mathbb{N}_0$ is **dual feasible** if, for every $W_S \subseteq W$:

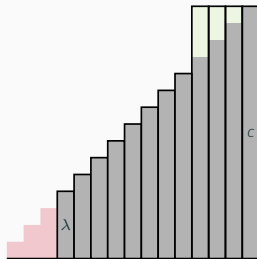
$$\sum_{w \in W_S} w \leq c \quad \Rightarrow \quad \sum_{w \in W_S} f(w) \leq f(c)$$

The combined lower bound L_f is:

$$L_f(c, W) = \left\lceil \frac{1}{f(c)} \sum_{w \in W} f(w) \right\rceil$$

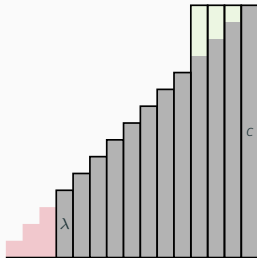






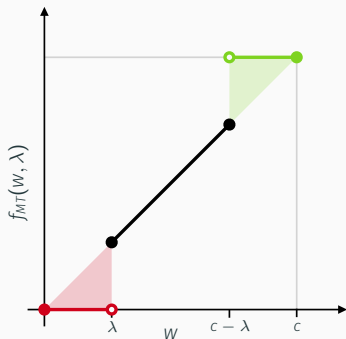
A basic DFF is f_{MT} depending on an integer parameter $0 \leq \lambda \leq \frac{c}{2}$:

$$f_{MT}(w, \lambda) = \begin{cases} 0 & \text{if } w < \lambda \\ c & \text{if } w > c - \lambda \\ w & \text{otherwise} \end{cases}$$



A basic DFF is f_{MT} depending on an integer parameter $0 \leq \lambda \leq \frac{c}{2}$:

$$f_{MT}(w, \lambda) = \begin{cases} 0 & \text{if } w < \lambda \\ c & \text{if } w > c - \lambda \\ w & \text{otherwise} \end{cases}$$

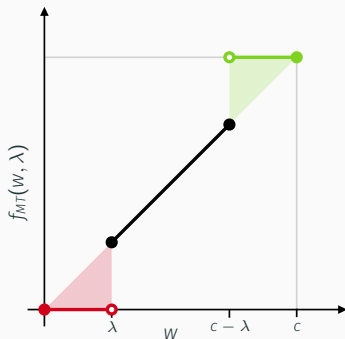


New classes of fast lower bounds for bin packing problems, Fekete et al., 2001 [🔗](#)

A basic DFF is f_{MT} depending on an integer parameter $0 \leq \lambda \leq \frac{c}{2}$:

$$f_{MT}(w, \lambda) = \begin{cases} 0 & \text{if } w < \lambda \\ c & \text{if } w > c - \lambda \\ w & \text{otherwise} \end{cases}$$

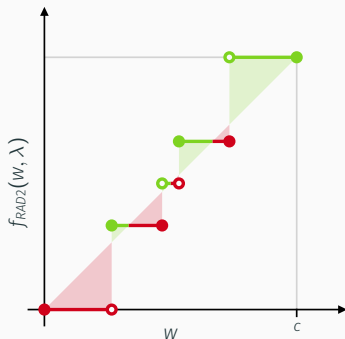
$$L_{MT}(c, W) = \max_{0 \leq \lambda \leq \frac{c}{2}} \left[\frac{1}{f_{MT}(c, \lambda)} \sum_{w \in W} f_{MT}(w, \lambda) \right]$$



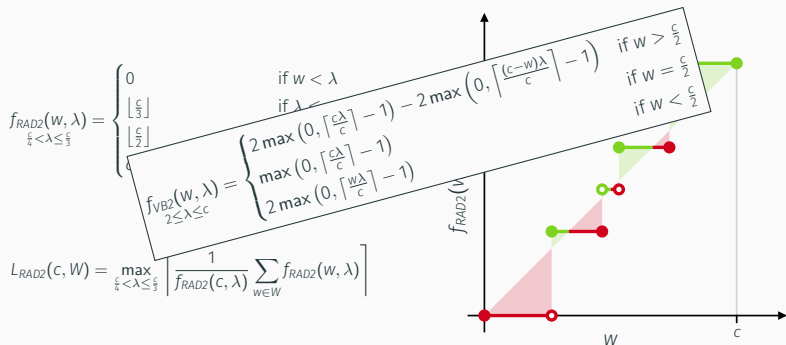
Another basic DFF is f_{RAD2} with parameter $\frac{c}{4} < \lambda \leq \frac{c}{3}$:

$$f_{RAD2}(w, \lambda) = \begin{cases} 0 & \text{if } w < \lambda \\ \lfloor \frac{c}{3} \rfloor & \text{if } \lambda \leq w \leq c - 2\lambda \\ \lfloor \frac{c}{2} \rfloor & \text{if } c - 2\lambda < w < 2\lambda \\ c - f_{RAD2}(c - w, \lambda) & \text{if } 2\lambda \leq w \end{cases}$$

$$L_{RAD2}(c, W) = \max_{\frac{c}{4} < \lambda \leq \frac{c}{3}} \left[\frac{1}{f_{RAD2}(c, \lambda)} \sum_{w \in W} f_{RAD2}(w, \lambda) \right]$$

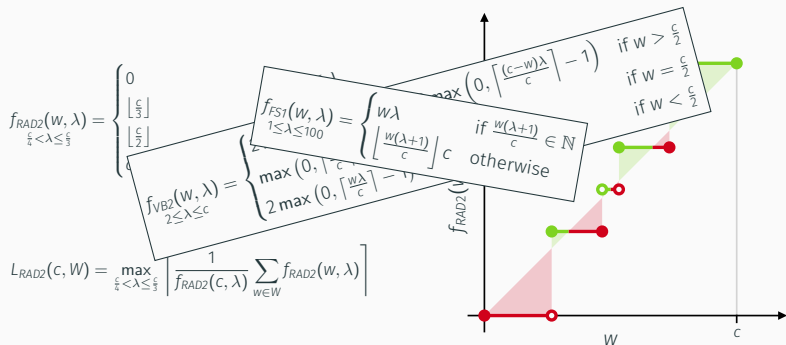


Another basic DFF is f_{RAD2} with parameter $\frac{c}{4} < \lambda \leq \frac{c}{3}$:



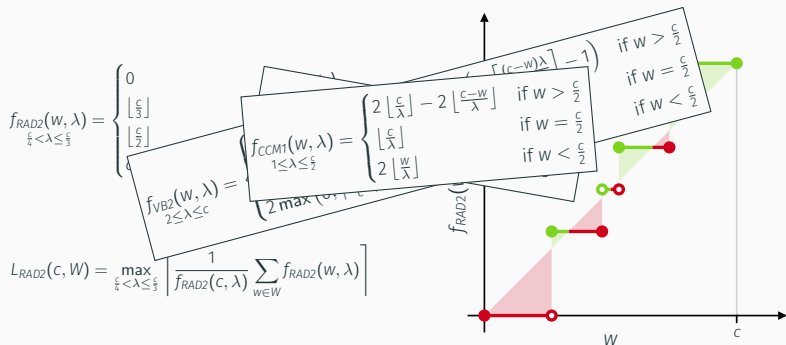
A survey of dual-feasible and superadditive functions, Clautiaux et al., 2010 [↗](#)

Another basic DFF is f_{RAD2} with parameter $\frac{c}{4} < \lambda \leq \frac{c}{3}$:



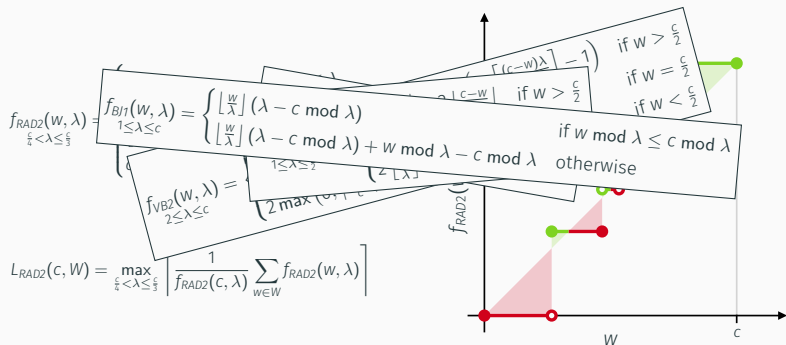
A survey of dual-feasible and superadditive functions, Clautiaux et al., 2010 [🔗](#)

Another basic DFF is f_{RAD2} with parameter $\frac{c}{4} < \lambda \leq \frac{c}{3}$:



A survey of dual-feasible and superadditive functions, Clautiaux et al., 2010 [🔗](#)

Another basic DFF is f_{RAD2} with parameter $\frac{c}{4} < \lambda \leq \frac{c}{3}$:



A survey of dual-feasible and superadditive functions, Clautiaux et al., 2010 [\[1\]](#)

Which DFF is the best?

Lower bound	Total Optimal	Only Optimal
L_{CCM1}	1219	40
L_{MT}	1151	2
L_{BJ1}	1101	47
L_{VB2}	973	1
L_{FS1}	742	2
L_{RAD2}	189	10

None dominate.

Which DFF is the best?

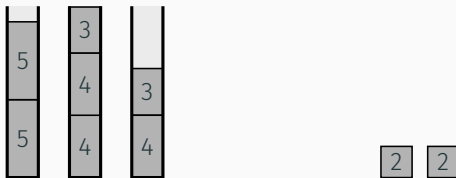
Lower bound	Total Optimal	Only Optimal
L_{CCM1}	1219	40
L_{MT}	1151	2
L_{BJ1}	1101	47
L_{VB2}	973	1
L_{FS1}	742	2
L_{RAD2}	189	10

Use the **GPU** to calculate **all** the lower bounds in **parallel!**

Lower bounds apply to instances, while search handles partial packings.

Reduce a partial packing to an “equivalent” instance.

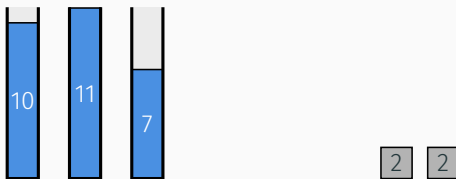
The basic reduction is R_0 , but R_{Min} and R_{Max} are generally more accurate.



Lower bounds apply to instances, while search handles partial packings.

Reduce a partial packing to an “equivalent” instance.

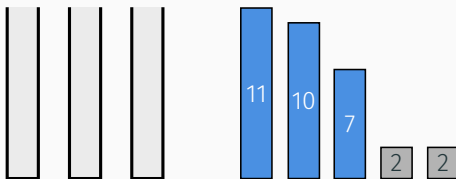
The basic reduction is R_0 , but R_{Min} and R_{Max} are generally more accurate.



Lower bounds apply to instances, while search handles partial packings.

Reduce a partial packing to an “equivalent” instance.

The basic reduction is R_0 , but R_{Min} and R_{Max} are generally more accurate.



Sequential calculation of the DFFs-based lower bound:

```

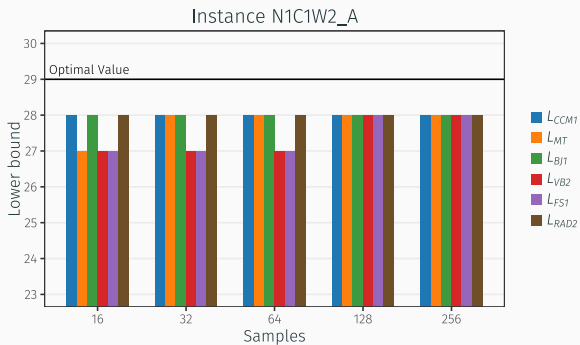
Function: getLowerBound(c, W, k, X) → lb
1 lb ← 0
2 for R ∈ {R0, RMin, RMax} do
3   (cR, WR) ← R(c, W, X)
4   for f ∈ {fCCM1, fMT, fBJ1, fVB2, fFS1, fRAD2} do
5     (λ̲, λ̄) ← getParametersMinMax(f, cR)
6     Lf ← 0
7     for λ ← λ̲ to λ̄ do
8       Lf ← max (Lf,  $\lceil \frac{1}{f(c, \lambda)} \sum_{w \in W} f(w, \lambda) \rceil$ )
9       lb ← max(lb, Lf)
10      if lb > k then return lb
11 return lb
    
```

Sequential calculation of the DFFs-based lower bound:

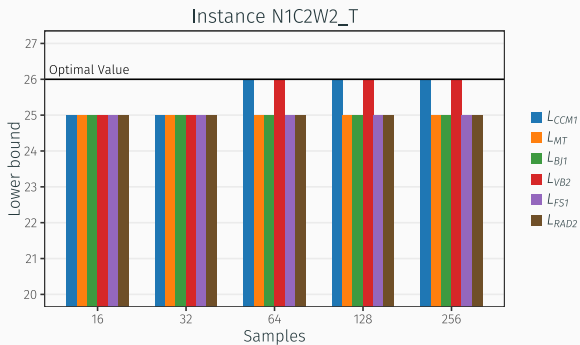
```
Function: getLowerBound( $c, W, k, X$ )  $\rightarrow lb$   
1  $lb \leftarrow 0$   
2 for  $R \in \{R_0, R_{Min}, R_{Max}\}$  do  
3    $(c_R, W_R) \leftarrow R(c, W, X)$   
4   for  $f \in \{f_{CCM1}, f_{MT}, f_{BJ1}, f_{VB2}, f_{FS1}, f_{RAD2}\}$  do  
5      $(\underline{\lambda}, \bar{\lambda}) \leftarrow getParametersMinMax(f, c_R)$   
6      $L_f \leftarrow 0$   
7     for  $\lambda \leftarrow \underline{\lambda}$  to  $\bar{\lambda}$  do  
8        $L_f \leftarrow \max(L_f, \lceil \frac{1}{f(c, \lambda)} \sum_{w \in W} f(w, \lambda) \rceil)$   
9      $lb \leftarrow \max(lb, L_f)$   
10    if  $lb > k$  then return  $lb$   
11 return  $lb$ 
```

- Independent iterations.
- Fits in L1 cache.
- No if/else branches.
- Few atomics.

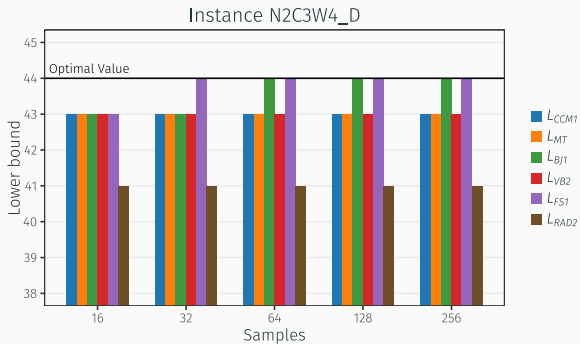
Lower bounds for different sampling of λ values:



Lower bounds for different sampling of λ values:



Lower bounds for different sampling of λ values:



Results

We used various BPP benchmarks from the literature to compare:

- Standard lower bound (L2).
- Linear relaxation of the Arc-Flow model (Arc-Flow)¹.
- Sampled sequential implementation (DFFs-Seq-CPU).
- Sampled multithread implementation (DFFs-Par-CPU).
- Complete GPU implementation (DFFs-GPU).

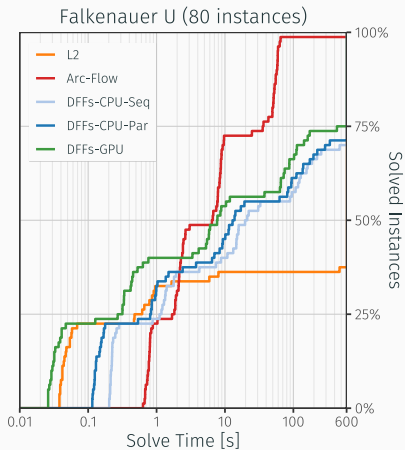
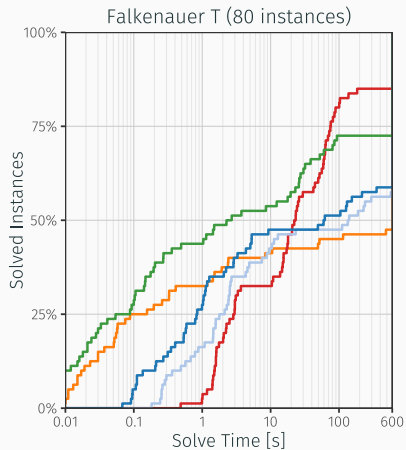
¹ Propagating the Bin Packing Constraint Using Linear Programming, Cambazard et al., 2010 [🔗](#)

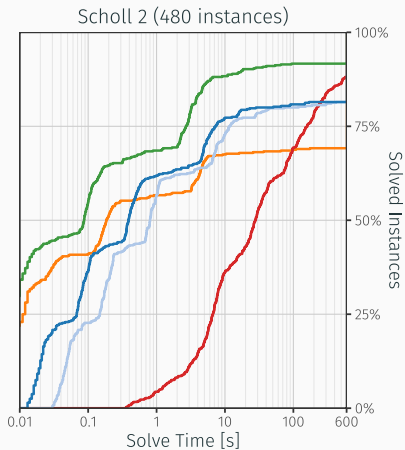
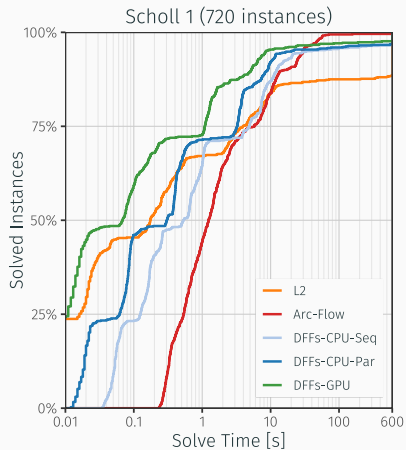
Model and search setup:

- Add one more bin until solution is found.
- Decreasing Best Fit heuristic.
- Symmetry breaking and dominance rules.
- 10 minutes timeout.

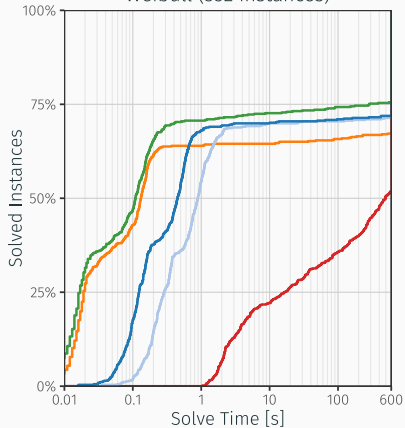
Benchmark system:

- Intel Core i7-10700K (8 Cores at 3.8 GHz).
- NVIDIA GeForce RTX 3080 (8704 CUDA Cores at 1.4 GHz).
- 32 GB RAM.
- Ubuntu Linux 22.04 LTS.
- CUDA 11.8.
- CPLEX 22.1 (for Arc-Flow).

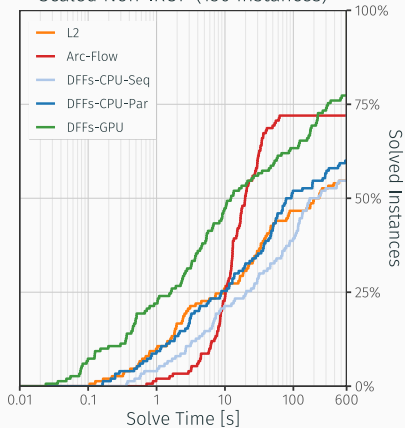




Weibull (552 instances)



Scaled Non-IRUP (150 instances)



Conclusion

Code repository:

<https://bitbucket.org/constraint-programming/bpp-minicpp>

<https://bitbucket.org/constraint-programming/fzn-minicpp>

Transparent GPU utilization in MiniZinc ():

```
include "globals.mzn";  
include "minicpp.mzn";  
  
int: n_items; % Number of items  
int: capacity; % Capacity of each bin  
array [1..n_items] of int: weights; % Weights of items  
...  
  
constraint bin_packing(capacity, items_bin, weights) ::gpu;  
...
```

BinPacking constraint:

- Novel trade-off pruning time vs pruning strength.
- Applicable to 2D, 3D, and Vector Packing Problems.

GPU in Constraint Programming:

- GPU can also be utilized for pruning.
- It is essential to **rethink algorithmic design** choices.