

Automatic Configuration of Benchmark Sets for Classical Planning

Álvaro Torralba,¹ Jendrik Seipp,² Silvan Sievers²

¹Aalborg University, Denmark

²University of Basel, Switzerland

October 21, 2020

Outline

- 1 The ICAPS Way
- 2 Benchmark Design Principles
- 3 Benchmark Configuration
- 4 Evaluation
- 5 Conclusion

The Cycle of Life (in Planning Research)

Everything you Always Wanted to Know About Planning
(But Were Afraid to Ask) — (Jörg Hoffmann, 2011)

The Cycle of Life (in Planning Research)

Everything you Always Wanted to Know About Planning (But Were Afraid to Ask) — (Jörg Hoffmann, 2011)

Is, then, the life of a researcher in heuristic search planning characterized by the following pseudo-code?

```
while ( not retired ) do  
    think up some new heuristic  $h^{foo-bar}$   
    run it on the benchmarks  
endwhile
```

Fig. 3. The life of a planning researcher?

The Cycle of Life (in Planning Research)

Everything you Always Wanted to Know About Planning (But Were Afraid to Ask) — (Jörg Hoffmann, 2011)

Is, then, the life of a researcher in heuristic search planning characterized by the following pseudo-code?

```
while ( not retired ) do  
    think up some new heuristic  $h^{foo-bar}$   
    run it on the benchmarks  
endwhile
```

Fig. 3. The life of a planning researcher?

The answer to that one is “NO!”. Far beyond just improving performance on benchmarks, the *understanding* of heuristics is where heuristic search planning really turns into a natural science. Dramatic progress has been made, in that science, during the last years. For example, Bonet and Geffner [3] proved

Empirical Evaluation – Examples from HSDIP'20

	Coverage			Time Score		
	offline	re-winds	on-stable	offline	re-winds	on-stable
agriculture (26)	0	0	0	0.0	0.0	0.0
airport (26)	34	34	34	2.5	18.5	26.1
barman (34)	4	0	4	0.3	0.0	2.0
blocks (36)	28	20	28	2.2	19.5	29.1
childtrack (26)	0	0	0	0.0	0.0	0.0
data-network (26)	14	11	14	1.1	7.2	12.1
depot (22)	13	6	13	1.0	3.3	9.6
driverlog (26)	15	7	15	1.1	4.3	10.3
elevators (26)	44	12	44	3.4	2.6	25.9
floortile (26)	6	0	6	0.3	0.0	0.8
freecell (26)	68	30	68	4.8	10.9	35.8
grid (26)	19	7	19	1.5	4.7	12.1
grid (5)	3	1	3	0.2	1.0	2.2
grripper (26)	8	6	8	0.6	4.8	7.0
hiking (26)	14	8	15	1.0	5.1	10.9
logistics (26)	39	19	39	2.8	11.8	25.0
microworld (26)	144	133	144	11.2	80.6	131.6
movie (26)	30	30	30	2.4	42.7	42.4
mpm (26)	29	24	29	2.3	19.1	26.5
mystery (26)	19	15	19	1.5	12.6	17.8
nonstop (26)	20	12	20	1.5	8.0	15.3
openstacks (26)	53	21	53	3.8	12.0	29.8
organic (26)	7	7	7	0.5	6.0	6.1
organic-split (26)	10	6	10	0.7	1.9	4.2
parcprinter (26)	38	34	38	2.9	29.2	34.1
parking (26)	13	1	13	0.9	0.1	5.4
pathways (26)	5	4	5	0.3	5.1	5.5
pegel (26)	48	42	48	3.7	22.8	35.6
petri-net (26)	0	0	0	0.0	0.0	0.0
pipes-04 (26)	25	14	25	1.8	10.4	19.0
pipes-4 (26)	18	8	18	1.3	5.1	12.1
psr-small (26)	50	40	50	3.9	48.2	54.8
rovers (26)	8	7	8	0.6	6.6	8.1
satellite (26)	7	6	7	0.5	5.5	7.3
scanalyzer (26)	35	7	35	2.7	5.7	19.6
snake (26)	12	6	12	0.9	2.5	7.4
sokoban (26)	50	33	50	3.8	19.6	39.5
spider (26)	15	7	15	1.1	2.9	8.6
stomp (26)	16	14	16	1.2	12.5	17.1
termex (26)	12	0	12	0.8	0.0	3.2
terris (26)	11	3	10	1.1	0.8	1.3
tillyost (26)	25	18	24	1.8	5.8	15.3
tpp (26)	8	7	8	0.6	8.0	8.9
transport (26)	34	20	36	2.6	10.4	22.4
trucks (26)	13	9	13	0.9	5.3	8.5
visual (26)	30	33	30	2.3	27.8	30.3
woodwork (26)	49	38	49	3.8	24.3	40.5
zonetravel (26)	13	7	13	1.0	4.9	8.7
Sum (2627)	1156	766	1155	1189	86.8	539.8

Empirical Evaluation – Examples from HSDIP'20

	Coverage			Time Score		
	offline	re-winds	on-stable	offline	re-winds	on-stable
agricola (26)	0	0	0	0.0	0.0	0.0
airport (26)	34	34	34	2.5	19.5	26.1
barman (34)	4	0	4	0.3	0.0	2.0
blocks (36)	28	20	28	2.2	19.5	29.3
chilkick (26)	0	0	0	0.0	0.0	0.0
data-network (26)	14	11	14	1.1	7.2	12.1
depot (22)	13	6	13	1.0	3.3	9.6
driveblog (26)	15	7	15	1.1	4.3	10.5
elevators (26)	44	12	44	3.4	2.6	25.9
floortile (26)	6	0	6	0.3	0.0	0.8
freecell (26)	68	30	68	4.8	10.9	35.8
grid (26)	19	7	19	1.5	4.7	12.2
grid (5)	3	1	3	0.2	1.0	2.2
grripper (26)	8	6	8	0.6	4.8	7.0
hiking (26)	14	8	15	1.0	5.1	10.9
logistics (26)	39	19	39	2.8	11.8	25.2
microworld (26)	144	133	144	11.2	80.6	131.6
movie (26)	30	30	30	2.4	42.7	42.4
mpm (26)	29	24	29	2.3	19.1	26.5
mystery (26)	19	15	19	1.5	12.6	17.8
nonstary (26)	20	12	20	1.5	8.0	15.3
openstacks (26)	53	2	53	3.8	12.0	29.8
organic (26)	7	7	7	0.5	6.0	6.1
organic-split (26)	10	6	10	0.7	1.9	4.2
parcprinter (26)	38	34	38	2.9	28.2	34.2
parking (26)	13	1	13	0.9	0.1	5.4
pathways (26)	5	4	5	0.3	5.1	5.5
pegel (26)	48	42	48	3.7	22.8	35.6
petri-net (26)	0	0	0	0.0	0.0	0.0
peps-4 (26)	25	14	25	1.8	10.4	19.0
peps-4 (26)	18	8	18	1.3	5.1	12.1
pet-small (26)	50	49	50	3.9	48.2	54.5
rovers (26)	8	7	8	0.6	6.6	8.1
santitas (26)	7	6	7	0.5	5.5	7.3
scanalyzer (26)	35	7	33	2.7	5.7	19.6
snake (26)	12	6	12	0.9	2.5	7.4
sokoban (26)	50	33	50	3.8	19.6	39.5
spider (26)	15	7	15	1.1	2.9	8.6
strong (26)	16	16	16	1.2	12.5	17.1
termis (26)	12	0	12	0.8	0.0	3.2
terris (26)	11	3	10	1.1	0.8	1.3
tillybot (26)	25	24	25	1.8	5.8	15.4
tip (26)	8	7	8	0.6	8.0	8.9
transport (26)	34	20	36	2.6	10.4	22.4
trucks (26)	15	9	15	0.9	5.3	8.5
visual (26)	30	33	30	2.3	27.8	30.3
woodwork (26)	49	38	49	3.8	24.3	40.5
worldtravel (26)	15	7	15	1.0	4.9	8.7
Sum (26*7)	1156	766	1155	1189	86.8	539.8

	Min	SPQR (100)			SPQR (100)		
		Q	R	S	Q	R	S
agricola	0	0	10	16	1	13	12
airport	23	35	36	33	1	21	13
barman	0	0	12	20	0	0	0
blocks	18	35	35	35	21	33	27
chilkick	0	0	1	0	0	0	0
data-network	0	1	5	2	0	0	0
depot	4	14	18	16	6	5	5
driveblog	7	19	18	18	8	12	13
elevators	0	5	38	7	0	0	0
floortile	0	0	2	1	0	0	0
freecell	20	46	79	80	15	80	23
grid	0	20	20	20	0	20	0
grripper	1	3	4	8	0	3	4
hiking	8	20	20	20	8	20	20
logistics	2	3	20	20	2	7	5
logistics	2	7	29	15	2	2	3
maintenace	0	14	11	14	0	0	0
microworld	55	150	150	150	71	150	150
movie	30	30	30	30	30	30	30
mpm	20	21	32	22	5	18	13
openstacks	15	15	17	14	0	6	12
openstacks	0	0	2	20	0	1	8
organic-synthesis	3	3	2	3	3	3	2
parcprinter	0	12	38	18	0	0	1
parking	0	0	7	0	0	0	0
pathways	4	5	10	8	4	4	4
pegel	17	20	20	18	20	20	20
petri-net	12	22	23	27	12	15	20
peps	40	40	40	40	40	40	40
peps	6	21	26	25	6	14	16
peps	6	15	27	12	7	15	8
pet-small	4	20	18	20	5	11	20
rovers	1	4	5	7	4	3	3
sokoban	6	13	19	10	8	11	12
spider	1	12	9	10	0	5	7
strong	18	19	19	17	15	19	19
termis	0	10	14	15	3	4	2
terris	0	20	9	20	2	11	19
tillybot	5	5	8	14	5	12	5
tip	3	19	16	20	0	1	7
transport	6	13	23	29	6	15	10
trucks	0	5	16	10	0	0	0
visual	6	9	15	9	6	6	7
visual	0	20	9	20	0	9	20
woodwork	1	1	2	4	1	1	1
worldtravel	8	20	20	20	7	9	9
Sum	359	775	933	965	346	651	536
>Qp					4	8	6
>R					2	5	4
>S					0	2	1

Empirical Evaluation – Examples from HSDIP'20

	Coverage			Time Score		
	offline	on-roads	on-stable	offline	on-roads	on-stable
agricola (20)	0	0	0	0.0	0.0	0.0
airport (20)	34	24	34	2.5	15.5	26.1
human (04)	4	0	4	0.3	0.0	2.0
blocks (70)	28	20	28	22	19.5	29.1
chickadee (20)	0	0	0	0.0	0.0	0.0
data-network (20)	14	11	14	1.1	7.2	12.1
depot (12)	13	6	13	1.0	3.3	9.6
drivlog (20)	15	7	15	1.1	4.3	10.5
elevators (04)	44	12	44	3.4	2.6	25.9
floortile (04)	6	0	6	0.3	0.0	0.8
freecell (04)	68	30	68	4.8	10.9	35.8
gfd (20)	19	7	19	1.5	4.7	12.2
grid (5)	3	1	3	0.2	1.0	2.2
grripper (20)	8	6	8	0.6	4.8	7.0
hiking (20)	14	8	15	1.0	5.1	10.7
logistics (04)	20	19	20	2.8	11.8	25.1
microworld (04)	144	133	144	11.2	80.6	131.6
novine (20)	30	20	30	2.4	42.7	42.4
expmpe (04)	29	24	29	2.3	19.1	26.5
novosty (04)	19	15	19	1.5	12.6	17.8
mystery (20)	20	12	20	1.5	8.0	15.3
openstacks (04)	53	22	53	3.8	12.0	29.8
organic (20)	7	7	7	0.5	6.0	6.1
organic-split (20)	10	6	10	0.7	1.9	4.2
parcprinter (04)	38	34	38	2.9	28.2	34.2
parking (04)	13	1	13	0.9	0.1	6.4
pathways (04)	5	4	5	0.3	5.1	5.5
pegel (04)	48	42	48	3.7	22.8	35.6
petri-net (20)	0	0	0	0.0	0.0	0.0
pipes-01 (20)	25	14	25	1.8	10.4	19.0
pipes-1 (20)	18	8	18	1.3	5.1	12.1
pre-small (04)	50	49	50	3.9	42.4	54.5
revers (04)	8	7	8	0.6	6.6	8.1
satlite (20)	7	6	7	0.5	3.5	7.2
scary (04)	3	3	3	2.7	5.7	19.6
scarylizer (20)	12	6	12	0.9	5.5	7.6
sokoban (20)	50	33	50	3.8	19.6	39.5
spider (20)	15	7	15	1.1	2.9	8.6
storage (04)	16	14	16	1.2	12.5	17.1
termes (20)	12	0	12	0.8	0.0	3.2
terris (07)	11	3	10	1.1	0.8	1.3
tidytot (04)	25	18	25	1.8	5.8	11.4
tip (04)	8	7	8	0.6	8.0	8.9
transport (20)	34	20	36	2.6	10.4	22.4
trucks (04)	13	9	13	0.9	5.3	8.5
visual (04)	20	33	30	2.7	27.8	30.3
woodwork (04)	49	38	49	3.8	24.3	40.3
zenotravel (20)	13	7	13	1.0	4.9	8.7
Sum (087)	1156	766	1155	1159	86.8	539.8

	Time		SPOG (Mbps)	
	min	max	min	max
agricola	0	0	0	0
airport	13	35	36	35
human	0	0	0	0
blocks	18	35	35	35
chickadee	0	0	0	0
data-network	0	1	5	2
depot	4	14	18	6
drivlog	7	19	18	8
elevators	0	5	20	7
floortile	0	0	2	4
freecell	20	46	79	80
gfd	0	20	20	20
grid	1	3	4	8
grripper	8	20	20	9
hiking	2	3	20	2
logistics	2	7	29	15
microworld	15	150	150	150
novine	30	30	30	30
expmpe	20	21	32	22
novosty	15	15	17	14
mystery	0	0	2	20
openstacks	0	3	2	3
organic-split	0	12	20	18
parcprinter	0	0	7	0
parking	4	5	10	8
pathways	17	20	20	18
pegel	12	22	25	27
petri-net	0	0	0	0
revers	6	21	26	25
satlite	6	15	27	12
scarylizer	4	20	20	5
scary	4	4	5	7
spider	1	12	9	19
storage	14	15	19	17
termes	0	0	10	14
terris	0	20	9	20
tidytot	5	5	8	14
tip	3	3	19	36
transport	6	5	13	25
trucks	6	9	15	9
visual	0	20	0	20
woodwork	1	1	2	4
zenotravel	8	20	20	37
SEM	359	775	933	945

Coverage	arb	inv	rnd	GZD	BD	ZCA	VDM	ZCP	AM	GZD+BD	Scorpion
Airport (30)	28	27	23	29	27	28	24	24	27	28	29
Blocks (15)	28	27	28	28	28	28	28	28	28	28	28
DataNetwork (20)	12	12	12	13	12	12	12	12	12	13	14
Depot (22)	7	7	7	7	7	7	7	7	7	10	13
DriverLog (20)	13	14	13	14	13	13	13	13	13	13	15
Elevators (30)	22	22	20	22	22	22	22	22	22	22	24
Freecell (04)	15	15	15	16	12	15	15	15	21	33	64
Grid (5)	2	2	1	2	2	2	2	2	2	2	3
Hiking (20)	10	10	9	10	9	8	9	9	10	9	14
Logistics (03)	27	27	25	27	25	25	25	25	25	27	34
Mprime (35)	23	25	22	23	23	23	22	22	25	24	31
Mystery (19)	16	17	15	17	17	17	17	17	17	17	19
Novosty (20)	16	16	14	17	15	14	18	18	16	18	20
Openstacks (08)	31	31	31	31	31	30	31	31	31	31	34
OrgSynth-split (20)	15	14	15	14	10	15	15	15	15	15	10
Parcprinter (30)	19	22	19	22	19	19	22	22	18	20	30
Parking (40)	9	9	6	10	10	10	12	12	8	13	13
Pegel (16)	35	34	33	35	34	34	35	34	34	35	35
Pipes-notank (50)	18	18	17	18	17	17	18	18	17	18	25
Pipes-tank (50)	12	10	12	11	11	9	12	12	12	12	18
PVCs-Alignment (20)	9	9	7	9	9	9	9	9	9	9	9
Revers (40)	9	11	9	9	9	9	9	9	9	9	9
Satellite (30)	8	12	7	8	14	13	15	15	10	14	8
Scarylizer (30)	16	16	16	16	15	14	16	16	16	16	18
Snake (20)	6	6	4	6	4	6	6	6	6	7	13
Sokoban (30)	30	29	30	30	30	30	30	30	30	30	30
Spider (20)	11	11	9	12	11	9	11	11	10	12	15
Termes (20)	7	6	6	7	6	6	7	6	7	6	13
Tidytot (40)	23	22	22	23	22	15	22	22	22	22	22
VisitAll (40)	16	15	17	15	36	36	36	36	36	36	30
Woodworking (30)	19	22	19	20	22	20	22	20	22	20	30
Zenotravel (20)	13	13	12	13	12	12	13	13	13	12	13
Others (601)	331	331	331	331	331	331	331	331	331	331	346
Sum (1672)	856	865	821	873	869	850	884	883	858	914	1020

Empirical Evaluation – Examples from HSDIP'20

Table showing Coverage and Time Score for various benchmarks. Columns include 'coverage' (offline, on-track, on-table) and 'time score' (offline, on-track, on-table).

Table showing heuristic usage for various benchmarks. Columns include 'total', 'LRT', 'LMT', 'SNGSL (Mbps)', and 'usage' (%, %, %, %).

Table comparing Blind Search and Duplicate Checking for benchmarks. Columns include 'Domain #', 'Blind Search Pruning', 'Duplicate Checking', and 'A* with LRT Pruning Duplicate Checking'.

Table showing Coverage and Scorpion for various benchmarks. Columns include 'Coverage' (arb, inv, rnd, GZD, BD, ZCA, VIDM, ZCP, AM, GZD+BD) and 'Scorpion'.

Empirical Evaluation – The ICAPS/IPC Way

The ICAPS/IPC Way

- Measure coverage
- Time limit 30 minutes
- Memory limit 2-8 GB
- Use the benchmarks from the International Planning Competition

Empirical Evaluation – The ICAPS/IPC Way

The ICAPS/IPC Way

- Measure coverage
- Time limit 30 minutes
- Memory limit 2-8 GB
- Use the benchmarks from the International Planning Competition

Having a standard evaluation setting is generally beneficial:

- Reproducibility
- Interpretability
- Avoids hand picking results

Empirical Evaluation – The ICAPS/IPC Way

The ICAPS/IPC Way

- Measure coverage
- Time limit 30 minutes
- Memory limit 2-8 GB
- Use the benchmarks from the International Planning Competition

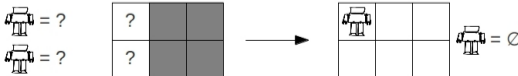
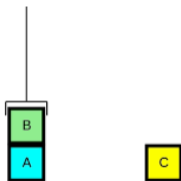
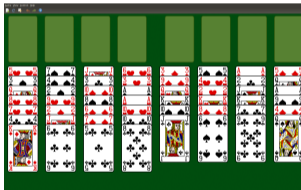
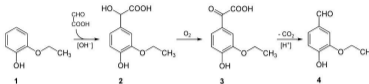
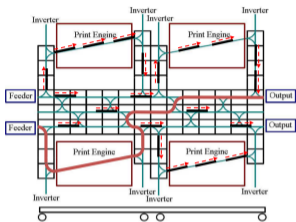
Having a standard evaluation setting is generally beneficial:

- Reproducibility
- Interpretability
- Avoids hand picking results

Outline

- 1 The ICAPS Way
- 2 Benchmark Design Principles**
- 3 Benchmark Configuration
- 4 Evaluation
- 5 Conclusion

The diversity in the IPC Benchmark Set



So, What's Wrong with the IPC Benchmark Set?

	IPC		
	L	D	O
Nomystery (20)	11	20	12
Rovers (40)	40	40	40
Woodworking (50)	50	50	50
Total	101	110	102

Table: Coverage of LAMA (L), Decstar (D) and OLCFF (O)

So, What's Wrong with the IPC Benchmark Set?

	IPC		
	L	D	O
Nomystery (20)	11	20	12
Rovers (40)	40	40	40
Woodworking (50)	50	50	50
Total	101	110	102

Table: Coverage of LAMA (L), Decstar (D) and OLCFF (O)

- Different number of instances per domain
- **Instance scaling**: too easy, too hard, and not smooth

So, What's Wrong with the IPC Benchmark Set?

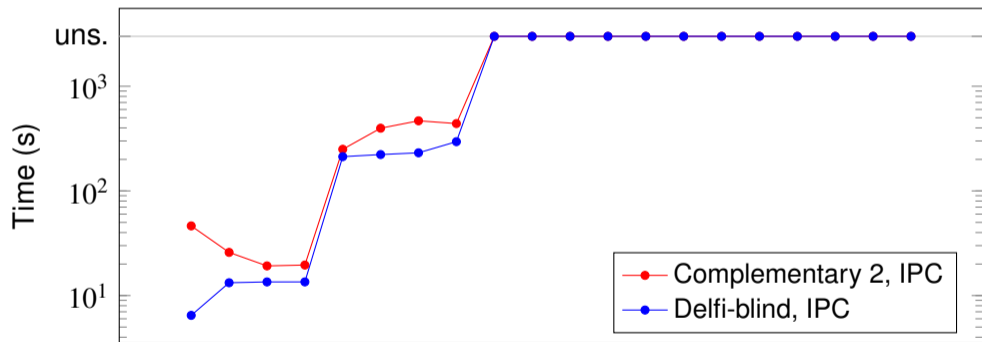
	IPC			New'14		
	L	D	O	L	D	O
Nomystery (20)	11	20	12	25	30	24
Rovers (40)	40	40	40	22	18	21
Woodworking (50)	50	50	50	18	27	30
Total	101	110	102	65	75	75

Table: Coverage of LAMA (L), Decstar (D) and OLCFF (O)

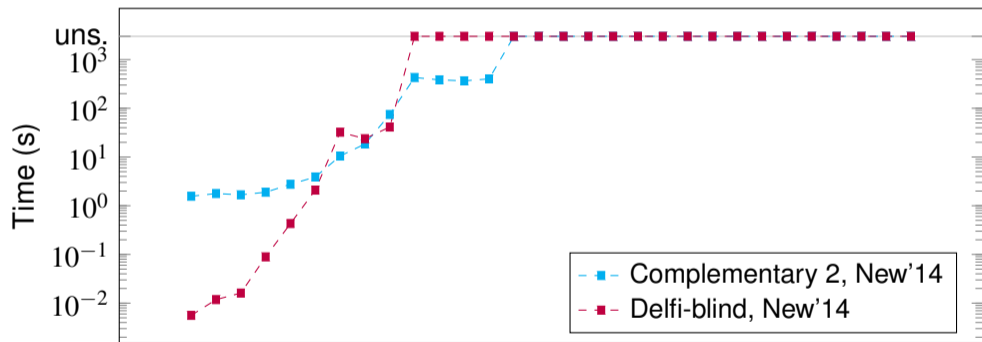
- Different number of instances per domain
- **Instance scaling**: too easy, too hard, and not smooth

→ Experiments on some domains of the IPC benchmark set may not observe any difference between planners even if it exists!

Non-Smooth Scaling



Smooth Scaling



Contribution

An **automatic tool** to select instances from a given domain (more informative than the IPC set to compare current and future planners)

Contribution

An **automatic tool** to select instances from a given domain (more informative than the IPC set to compare current and future planners)

- ① **Smooth** scaling from **easy** to **hard** instances:

Contribution

An **automatic tool** to select instances from a given domain (more informative than the IPC set to compare current and future planners)

- ① **Smooth** scaling from **easy** to **hard** instances:
 - Easy: solvable by any planner that anyone would compare against (**baseline**)
 - Hard: out of reach of **current existing planners** within a reasonable time limit

Contribution

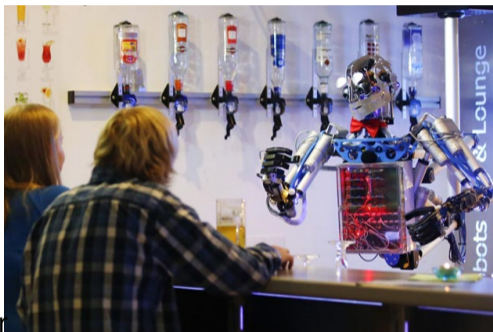
An **automatic tool** to select instances from a given domain (more informative than the IPC set to compare current and future planners)

- 1 **Smooth** scaling from **easy** to **hard** instances:
 - Easy: solvable by any planner that anyone would compare against (**baseline**)
 - Hard: out of reach of **current existing planners** within a reasonable time limit
- 2 Minimize bias towards/against planners used

Outline

- 1 The ICAPS Way
- 2 Benchmark Design Principles
- 3 Benchmark Configuration**
- 4 Evaluation
- 5 Conclusion

Example Domain: Barman



Instance Generator

```
./barman-generator.py <num_cocktails> <num_ingredients>  
                    <num_shots> [<random_seed>]  
  
num_cocktails (min 1)  
num_ingredients (min 2)  
num_shots (min num_cocktails+1)
```

Instance Generation Problem

Input:

- domain
- instance generator
- a baseline planner
- a set of state-of-the-art planners

Output: set of instances with a good scaling

Instance Generation Problem

Input:

- domain
- instance generator
- a baseline planner
- a set of state-of-the-art planners

Output: set of instances with a good scaling

Generate instances → Compute/Estimate runtimes → Select instances

Instance Generation Problem

Input:

- domain
- instance generator
- a baseline planner
- a set of state-of-the-art planners

Output: set of instances with a good scaling

Generate instances → Compute/Estimate runtimes → Select instances

How to avoid **bias** wrt. the set of considered planners?

Instance Generation Problem

Input:

- domain
- instance generator
- a baseline planner
- a set of state-of-the-art planners

Output: set of instances with a good scaling

Generate instances → Compute/Estimate runtimes → Select instances

How to avoid **bias** wrt. the set of considered planners?

Output: set of **linear scaling of parameters for the generator** that produce a good scaling in runtime

Sequences of instances

User specifies characteristics of the generator parameters:

Linear Attributes

cocktails	shots	ingredients
$b \in [1, 6]$	$b \in [1, 5]$	$v \in \{3, 4, 5\}$
$m \in [1, 5]$	$m \in [0, 5]$	
	+ cocktails	

Enumerated attributes

Sequences of instances

User specifies characteristics of the generator parameters:

Linear Attributes:

- Numeric value
- Increase size of the task
- User specifies ranges for the base value (b) and slope (m)

cocktails	shots	ingredients
$b \in [1, 6]$	$b \in [1, 5]$	$v \in \{3, 4, 5\}$
$m \in [1, 5]$	$m \in [0, 5]$ + cocktails	

Enumerated attributes

Sequences of instances

User specifies characteristics of the generator parameters:

Linear Attributes:

- Numeric value
- Increase size of the task
- User specifies ranges for the base value (b) and slope (m)

cocktails	shots	ingredients
$b \in [1, 6]$	$b \in [1, 5]$	$v \in \{3, 4, 5\}$
$m \in [1, 5]$	$m \in [0, 5]$ + cocktails	

Enumerated attributes:

- Finite set of values
- Fixed in the sequence

Sequences of instances

User specifies characteristics of the generator parameters:

Linear Attributes:

- Numeric value
- Increase size of the task
- User specifies ranges for the base value (b) and slope (m)

Enumerated attributes:

- Finite set of values
- Fixed in the sequence

cocktails	shots	ingredients
$b \in [1, 6]$	$b \in [1, 5]$	$v \in \{3, 4, 5\}$
$m \in [1, 5]$	$m \in [0, 5]$ + cocktails	
Our system may select sequences like:		
$(b = 5,$	$(b = 1, m = 0,$	$(v = 3)$
$m = 1.34)$	+cocktails)	
5	6	3
6	7	3
7	8	3
9	10	3
10	11	3
11	12	3
13	14	3

Optimization Process

- 1 Generate candidate sequences that scale smoothly
- 2 Choose selected (sub-)sequences to include easy to hard instances

Sequence Optimization

We use SMAC to optimize the value of b , m and v for each parameter

- Measure instance difficulty as the best runtime by any planner (time limit 180 seconds)
- Penalty based on smoothness of scaling difficulty (ideally by a factor of 1.5 to 2)

Runtimes: 10.36, 15.41, 18.9, 28.02, 29.27, 68.01

Ratios: 1.48, 1.22, 1.48, 1.04, 2.32

Penalties: 0.02, 0.54, 0.02, 0.91, 0.13

Total penalty 1.62

Sequence Selection

MIP encoding to select sequences satisfying **hard constraints**:

- There are 30 instances
- **(Easy)** Baseline solves at least one instance in less than 30 seconds
- **(Hard)** Sub-sequences go from easy ($\leq 180s$) to hard ($> 2000s$)
- **(Diverse)** Don't repeat the same parameters more than twice

and **soft constraints**:

- **(Easy)** Baseline solves 2 to 6 instances under 30 seconds
- **(Easy)** State-art planners solve 8 to 15 instances under 180 seconds
- **(Hard)** All sequences end in a very hard instance
- **(Diverse)** Don't repeat the same parameters more than once
- **(Smooth)** Minimize penalty of selected sequences

Outline

- 1 The ICAPS Way
- 2 Benchmark Design Principles
- 3 Benchmark Configuration
- 4 Evaluation
- 5 Conclusion

Experiments

Compare our new benchmark sets against the IPC

- 26 domains
- Satisficing and Optimal track
- 2 new benchmark sets that differ on the “training set”:
 - New’14: using planners up to 2014
 - New’20: using all available planners
- Evaluation based on planners from IPC’18

Evaluation Criteria

How to evaluate the quality of a benchmark set?

Evaluation Criteria

How to evaluate the quality of a benchmark set?

- Coverage range: generally better if all planners solve some instance and no planner solves all instances
- Comparisons: number of pairs (X, Y) of planners, such that $coverage(X) \neq coverage(Y)$

Evaluation Criteria

How to evaluate the quality of a benchmark set?

- Coverage range: generally better if all planners solve some instance and no planner solves all instances
- Comparisons: number of pairs (X, Y) of planners, such that $coverage(X) \neq coverage(Y)$

Goodhart's law: "When a measure becomes a target, it ceases to be a good measure." – Marilyn Strathern

→ Comparisons is a useful metric to compare benchmarks but not a metric to optimize for (would introduce **bias** towards the set of planners)

Results

Optimal	#IPC	coverage range			comparisons			Satisficing	#IPC	coverage range			comparisons		
		IPC	'14	'20	IPC	'14	'20			IPC	'14	'20	IPC	'14	'20
barman	34	4-11	9-13	9-12	12	21	19	barman	40	39-40	7-25	9-30	7	24	27
blocksworld	35	18-30	5-12	5-12	18	24	24	blocksworld	35	35-35	7-24	4-22	0	27	28
childsnaek	20	0-6	9-20	6-21	12	18	22	childsnaek	20	1-20	14-30	2-19	27	25	28
data-network	20	6-14	5-12	5-16	27	25	27	data-network	20	9-19	10-30	13-30	24	27	25
depot	22	5-14	9-25	8-16	26	26	24	depot	22	21-22	12-20	11-26	7	27	22
driverlog	20	7-15	6-30	5-18	22	26	25	driverlog	20	20-20	29-30	9-19	0	12	24
elevators	50	28-44	7-14	10-18	26	26	23	elevators	50	49-50	30-30	30-30	7	0	0
floortile	40	16-34	9-18	8-17	21	21	22	floortile	40	4-40	1-12	1-11	17	25	24
grid	5	1-3	6-26	4-21	19	28	27	grid	5	5-5	4-20	9-21	0	26	24
gripper	20	8-20	11-30	11-30	7	7	7	gripper	20	20-20	26-30	26-30	0	7	7
hiking	20	12-18	7-9	5-16	23	15	25	hiking	20	10-20	2-22	3-26	24	28	27
logistics	63	13-34	5-17	5-14	27	27	25	logistics	63	51-63	5-30	5-26	17	27	26
miconic	150	56-142	4-28	3-30	25	27	28	miconic	150	150-150	30-30	30-30	0	0	0
nomystery	20	8-20	3-27	5-21	18	28	27	nomystery	20	12-20	19-30	2-30	23	18	26
openstacks	130	42-71	4-11	3-7	24	18	7	openstacks	160	99-160	12-21	14-23	21	27	25
parking	40	0-15	11-18	12-21	28	24	23	parking	40	36-40	14-20	13-16	7	24	21
rovers	40	6-13	4-26	6-19	25	22	7	rovers	40	38-40	10-22	6-30	7	26	27
satellite	36	7-14	8-30	4-27	22	25	26	satellite	36	26-36	5-30	6-14	23	17	23
scanalyzer	50	21-33	6-16	7-15	27	24	24	scanalyzer	50	48-50	9-16	13-14	12	21	12
snake	20	7-14	5-20	7-19	22	24	21	snake	20	3-17	6-30	5-14	27	28	26
storage	30	15-18	9-25	2-19	21	27	26	storage	30	21-30	6-26	7-17	26	27	26
tpp	30	7-20	7-30	2-7	24	24	21	tpp	30	29-30	10-26	6-21	15	27	27
transport	70	24-35	5-30	8-19	21	18	22	transport	70	65-70	22-30	15-23	7	24	26
visitall	40	12-30	6-21	5-20	27	27	27	visitall	40	36-40	4-30	4-29	7	24	26
woodworking	50	38-50	16-25	10-14	22	26	24	woodworking	50	28-50	6-30	5-30	13	27	27
zenotravel	20	8-13	6-30	3-13	23	26	28	zenotravel	20	20-20	6-29	5-17	0	23	25

Highlight: SAT track

Satisficing	comparisons		
	IPC	'14	'20
gripper	0	7	7
miconic	0	0	0
elevators	7	0	0
blocksworld	0	27	28
driverlog	0	12	24
grid	0	26	24
zenotravel	0	23	25
barman	7	24	27
depot	7	27	22
parking	7	24	21
rovers	7	26	27
transport	7	24	26
visitall	7	24	26

Outline

- 1 The ICAPS Way
- 2 Benchmark Design Principles
- 3 Benchmark Configuration
- 4 Evaluation
- 5 Conclusion

Conclusion

- New tool to automatically select instances
 - Our tool consistently generates well-scaled instance sets that are **useful to evaluate current planners**
- New benchmark set significantly better than the IPC benchmark set, specially in the SAT/AGL track

Conclusion

- New tool to automatically select instances
 - Our tool consistently generates well-scaled instance sets that are **useful to evaluate current planners**
- New benchmark set significantly better than the IPC benchmark set, specially in the SAT/AGL track
- We need your feedback!
 - Do you find the results of our tool useful?
 - Is there any reason to prefer the IPC set over our new one?
 - Are there any constraints that we should take into account (in general or for specific domains)?