

# Preprocessing and Inprocessing Techniques in SAT

Armin Biere

Institute for Formal Models and Verification  
Johannes Kepler University  
Linz, Austria

joint work with Marijn Heule and Matti Järvisalo

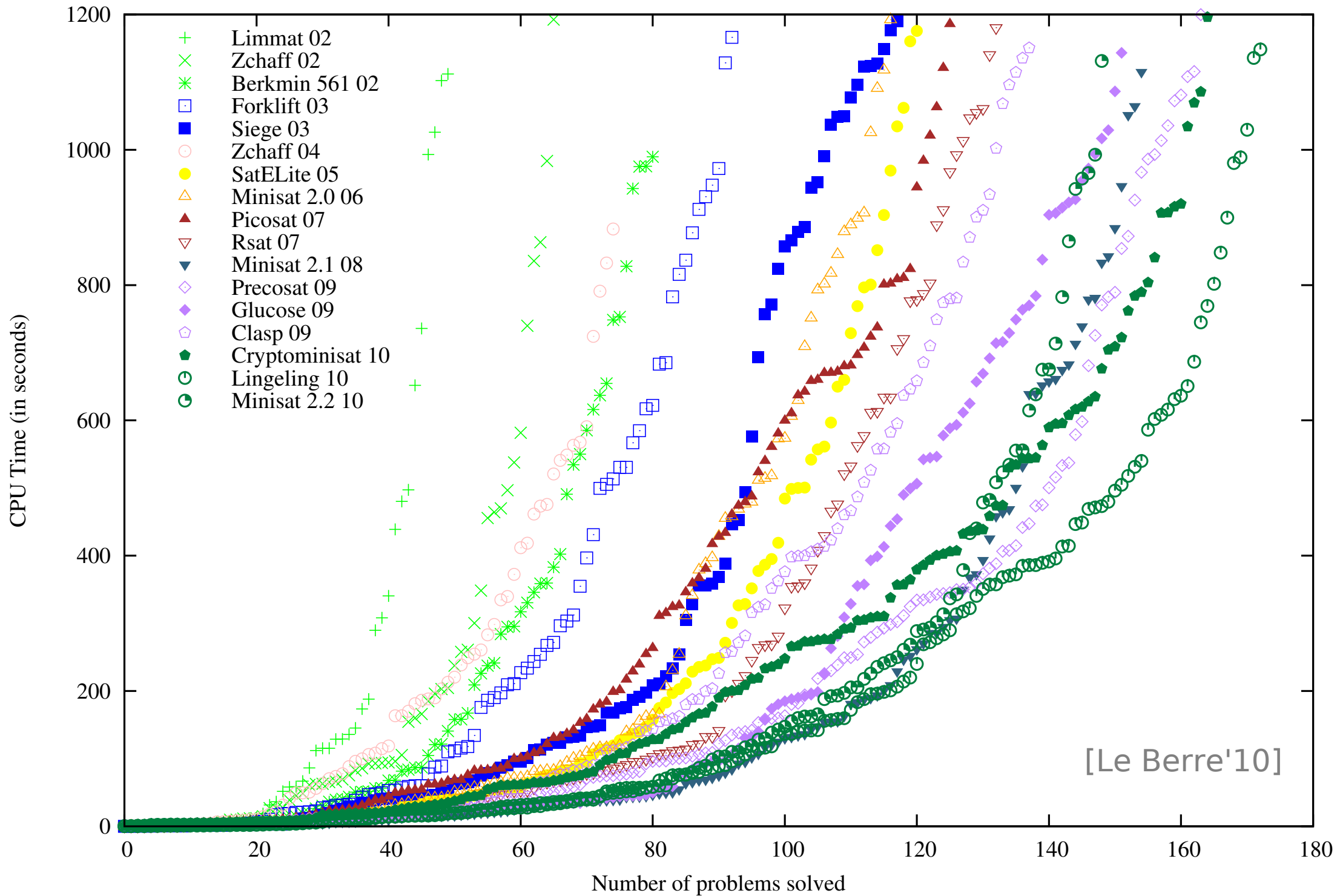
WorKer'11

3rd Workshop on Kernelization  
TU Vienna, Austria

Friday, September 2, 2011

- propositional logic:
  - variables **tie** **shirt**
  - negation  $\neg$  (not)
  - disjunction  $\vee$  disjunction (or)
  - conjunction  $\wedge$  conjunction (and)
- three conditions / clauses:
  - clearly one should not wear a **tie** without a **shirt**  $\neg\mathbf{tie} \vee \mathbf{shirt}$
  - not wearing a **tie** nor a **shirt** is impolite  $\mathbf{tie} \vee \mathbf{shirt}$
  - wearing a **tie** and a **shirt** is overkill  $\neg(\mathbf{tie} \wedge \mathbf{shirt}) \equiv \neg\mathbf{tie} \vee \neg\mathbf{shirt}$
- is the formula  $(\neg\mathbf{tie} \vee \mathbf{shirt}) \wedge (\mathbf{tie} \vee \mathbf{shirt}) \wedge (\neg\mathbf{tie} \vee \neg\mathbf{shirt})$  satisfiable?

Results of the SAT competition/race winners on the SAT 2009 application benchmarks, 20mn timeout



- failed literal *probing*
- variable elimination (VE)
- inprocessing
- lazy hyper binary resolution
- blocked clause elimination (BCE)
- hidden tautologies elimination (HTE)
- un hiding

*we are still working on tracking down the origin before* [Freeman'95] [LeBerre'01]

- key technique in look-ahead solvers such as Satz, OKSolver, March
  - failed literal probing at all search nodes
  - used to find the best decision variable and phase
- simple algorithm
  1. assume literal  $l$ , propagate (BCP), if this results in conflict, add unit clause  $\neg l$
  2. continue with all literals  $l$  until *saturation* (nothing changes)
- quadratic to cubic complexity
  - BCP linear in the size of the formula 1st linear factor
  - each variable needs to be tried 2nd linear factor
  - and tried again if some unit has been derived 3rd linear factor

- lifting
  - complete case split: literals implied in all cases become units
  - similar to Stålmarm's method and Recursive Learning [PradhamKunz'94]
- asymmetric branching
  - assume all but one literal of a clause to be false
  - if BCP leads to conflict remove originally remaining unassigned literal
  - implemented for a long time in MiniSAT but switched off by default
- generalizations:
  - vivification [PietteHamadiSais ECAI'08]
  - distillation [JinSomenzi'05][HanSomenzi DAC'07] probably most general (+ tries)

[Biere'04][SubbarayanPradhan'04][EénBiere SAT'05]

- goes back to original Davis & Putnam algorithm [DP'60]
  - eliminate variable  $x$  by adding all resolvents with  $x$  as pivot ...
  - ... and removing all clauses in which  $x$  or  $\neg x$  occurs
  - eliminating one variable is in the worst case quadratic
- bounded = apply only if increment in size is small
  - Quantor [Biere'03,Biere'04] bound increase in terms of literals (priority queue)
  - NiVER [SubbarayanPradhan'04] do not increase number of clauses (round-robin)
  - SatELite [EénBiere'05] do not increase number of clauses (priority queue)

- fast subsumption and strengthening [Biere'04][EénBiere'05]
  - backward subsumption: traverse clauses of least occurring literal
  - forward subsumption: one-watched literal scheme [Zhang'05]
  - 1st and 2nd level signatures = Bloom-filters for faster checking
  - strengthen clauses through self-subsuming resolution (later again)
  
- functional substitution
  - if  $x$  has a functional dependency, e.g. Tseitin translation of a gate
  - then only resolvents using exactly one “gate clause” need to be added

$$\overbrace{(\bar{x} \vee a)(\bar{x} \vee b)(x \vee \bar{a} \vee \bar{b})}^{x=a \wedge b} (x \vee c)(x \vee d)(\bar{x} \vee e)(\bar{x} \vee f) \quad 7 \text{ clauses}$$

$$(a \vee c)(a \vee d)(b \vee c)(b \vee d)(\bar{a} \vee \bar{b} \vee e)(\bar{a} \vee \bar{b} \vee f)(c \vee e)(c \vee f)(d \vee e)(d \vee f) \quad 6 + 4 \text{ clauses}$$



- **preprocessing** can be extremely beneficial
  - most SAT competition solvers use variable elimination (VE)  
[EénBiere SAT'05]
  - equivalence & XOR reasoning beneficial
  - probing / failed literal preprocessing / hyper binary resolution useful
  - however, even though polynomial, can not be run until completion
- **inprocessing**: simple idea to benefit from full preprocessing without penalty
  - “preempt” preprocessors after some time
  - resume preprocessing between restarts
  - limit preprocessing time in relation to search time

- allows to use costly preprocessors
  - without increasing run-time “much” in the worst-case
  - still useful for benchmarks where these costly techniques help
  - good examples: probing and distillation even VE can be costly
- additional benefit:
  - makes units / equivalences learned in search available to preprocessing
  - particularly interesting if preprocessing simulates encoding optimizations
- danger of hiding “bad” implementation though ...
- ... and hard(er) to debug

- one Hyper Binary Resolution step

[Bacchus-AAAI02]

$$\frac{(l \vee l_1 \vee \dots \vee l_n) \quad (\bar{l}_1 \vee l') \quad \dots \quad (\bar{l}_n \vee l')}{(l \vee l')}$$

- combines multiple resolution steps into one
- special case “hyper unary resolution” where  $l = l'$
- **HBR stronger than unit propagation** if it is repeated until (confluent) closure
- equality reduction: if  $(a \vee \bar{b}), (\bar{a} \vee b) \in f$  then replace  $a$  by  $b$  in  $f$

- can be simulated by unit propagation

[BacchusWinter-SAT03]

if  $(l \vee l') \in \text{HypBinRes}(f)$  then  $l' \in \text{UnitProp}(f \wedge \bar{l})$  or vice versa

- repeated probing, c.f. HypBinResFast

[GershmanStrichman-SAT05]

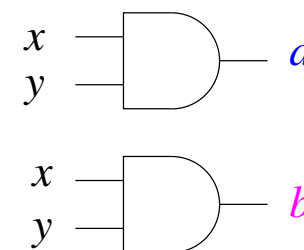
[BacchusWinter-SAT03][GershmanStrichman-SAT05]

- maintain acyclic and transitively-reduced binary implication graph
  - acyclic: decomposition in strongly connected components (SCCs)  
$$(\bar{a} \vee b)(\bar{b} \vee c)(\bar{c} \vee a) \wedge R \quad \text{equisatisfiable to} \quad R[a/b, a/c]$$
  - transitively-reduced: remove resp. do not add transitive edges
- not all literals have to be probed
  - if  $l \in \text{UnitProp}(r)$  and  $\text{UnitProp}(r)$  does not produce anything  
 $\Rightarrow$  no need to probe  $l$
  - at least with respect to units it is possible to focus on roots
- current algorithms too expensive to run until completion

- **time** complexity: seems to be at least quadratic, unfortunately also in practice
- **space** complexity: unclear, at most quadratic, linear?
- hyper binary resolution **simulates structural hashing** for AND gates  $a$  and  $b$

$$F \equiv (\bar{a} \vee x)(\bar{a} \vee y)(a \vee \bar{x} \vee \bar{y}) \quad (\bar{b} \vee x)(\bar{b} \vee y)(b \vee \bar{x} \vee \bar{y}) \quad \dots$$

$$\frac{(\bar{a} \vee x)(\bar{a} \vee y)(b \vee \bar{x} \vee \bar{y})}{(\bar{a} \vee b)} \quad \frac{(\bar{b} \vee x)(\bar{b} \vee y)(a \vee \bar{x} \vee \bar{y})}{(\bar{b} \vee a)}$$

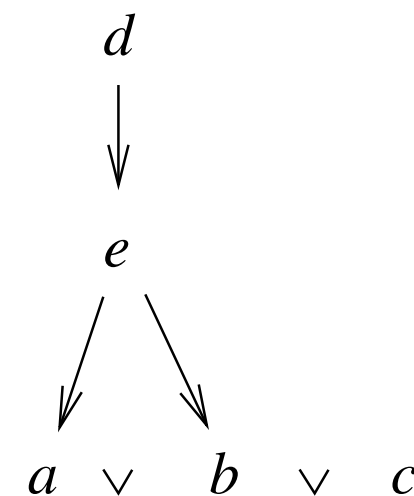


can also be seen by  $b \in \text{UnitProp}(F \wedge a)$  and  $a \in \text{UnitProp}(F \wedge b)$

- can not simulate structural hashing of XOR or ITE gates

- learn binary clauses lazily or on-the-fly
  - in the innermost (!) BCP loop
  - during BCP in search or during BCP in preprocessing (failed literal probing)
- whenever a **large** clause  $(a_1 \vee \dots \vee a_m \vee c)$  with  $m \geq 2$  becomes a reason for  $c$ 
  - partial assignment  $\sigma$  with  $\sigma(a_i) = 0$  and  $\sigma(c) = 1$
  - check whether exists literal  $d$  dominating all  $\bar{a}_i$
  - in implication graph restricted to binary clauses
  - which is a **tree** !
- learn  $(\bar{d} \vee c)$  if such a dominator exists

better  $(\bar{e} \vee c)$



1. trail contains assigned literals
2. set  $n_2$  and  $n_3$  to the trail level of those literals that still need to be propagated
3. while  $0 \leq n_3 \leq n_2 < |\text{trail}|$  and there is no conflict
  - (a) if  $n_2 < |\text{trail}|$ 
    - i. pick literal  $l$  at position  $n_2$ , increment  $n_2$  and visit binary clauses with  $\bar{l}$
    - ii. assign literals forced through these binary clauses first
  - (b) otherwise (necessarily  $n_3 < |\text{trail}|$ )
    - i. pick literal  $l$  at position  $n_3$ , increment  $n_3$  and visit large clauses with  $\bar{l}$
    - ii. assign literals forced through these large clauses

- for each assigned literal  $l$  keep **one** dominator  $\text{bindom}(l)$   
(in the implication graph restricted to binary clauses)
- thus  $\text{bindom}(l)$  is the root of binary implication tree of  $l$
- decisions  $l$  set  $\text{bindom}(l) = l$
- binary implications  $(a_1 \vee c)$  with  $\sigma(a_1) = 0, \sigma(c) = 1$  set  $\text{bindom}(c) = \text{bindom}(\overline{a_1})$
- necessary & sufficient for  $\overline{a_i}$  in large ( $m \geq 2$ ) reasons to have common dominator:  

$$(a_1 \vee \dots \vee a_m \vee c) \quad \text{bindom}(\overline{a_1}) = \dots = \text{bindom}(\overline{a_m})$$
- if condition triggers, actually use **least common ancestor** (closest dominator)
- use  $(\overline{d} \vee c)$  as **new reason** instead of  $(a_1 \vee \dots \vee a_m \vee c)$



- interleave search and preprocessing
  - bound time spent in search to roughly 80%
  - measured in number of propagations / resolutions
- during preprocessing / simplification on the top level
  - unit propagation on the top-level does LHBR top-level (1)
  - failed literal probing learns most binary clauses with LHBR probing (2)
- BCP during search learns binary clauses with LHBR search (3)

- rerunning SAT'09 competition with competition version 236 of PrecoSAT
  - 900 seconds time out
  - roughly twice as fast machines
- PrecoSAT without LHBR solves 6 less instances
  - 171 instead of 177 out of 292
- statistics
  - LHBR learned 48 million binary clause
  - on 292 instances that is 181 thousand learned binary clauses on average
  - additionally 202 million learned clauses through conflict analysis
  - 19% learned (binary) clauses due to LHBR

- no measurable overhead doing LHBR during BCP
  - so at least not harmful in contrast to many other “optimizations”
  - implementation in Lingeling became non-trivial
  - unfortunately does not simulate structural hashing in practice (!!)
- *not formally published* but implemented in PrecoSAT and Lingeling
  - source code of PrecoSAT available under MIT license
  - source code of Lingeling available under GPL license
  - extensions in [HanSomenzi-DATE'11]
- performs a limited version of on-the-fly strengthening / subsumption  
thus partially simulates distillation / vivification

blocked clause  $C \in F$

all clauses in  $F$  with  $\bar{l}$

fix a CNF  $F$

$$\bar{l} \vee \bar{a} \vee c$$

$$a \vee b \vee l$$

$$\bar{l} \vee \bar{b} \vee d$$

since all resolvents of  $C$  on  $l$  are tautological  $C$  can be removed

## Proof

assignment  $\sigma$  satisfying  $F \setminus C$  but not  $C$

can be extended to a satisfying assignment of  $F$  by flipping value of  $l$

**Definition** A literal  $l$  in a clause  $C$  of a CNF  $F$  **blocks**  $C$  w.r.t.  $F$  if for every clause  $C' \in F$  with  $\bar{l} \in C'$ , the resolvent  $(C \setminus \{l\}) \cup (C' \setminus \{\bar{l}\})$  obtained from resolving  $C$  and  $C'$  on  $l$  is a tautology.

**Definition** [Blocked Clause] A clause is **blocked** if has a literal that blocks it.

**Definition** [Blocked Literal] A literal is **blocked** if it blocks a clause.

**Example**

$$(a \vee b) \wedge (a \vee \bar{b} \vee \bar{c}) \wedge (\bar{a} \vee c)$$

only **first clause** is not blocked.

**second clause** contains two blocked literals:  $a$  and  $\bar{c}$ .

literal  $c$  in the **last clause** is blocked.

after removing either  $(a \vee \bar{b} \vee \bar{c})$  or  $(\bar{a} \vee c)$ , the clause  $(a \vee b)$  becomes blocked

actually all clauses can be removed

[JärvisaloBiereHeule-TACAS10]

**COI** Cone-of-Influence reduction

**MIR** Monotone-Input-Reduction

**NSI** Non-Shared Inputs reduction

---

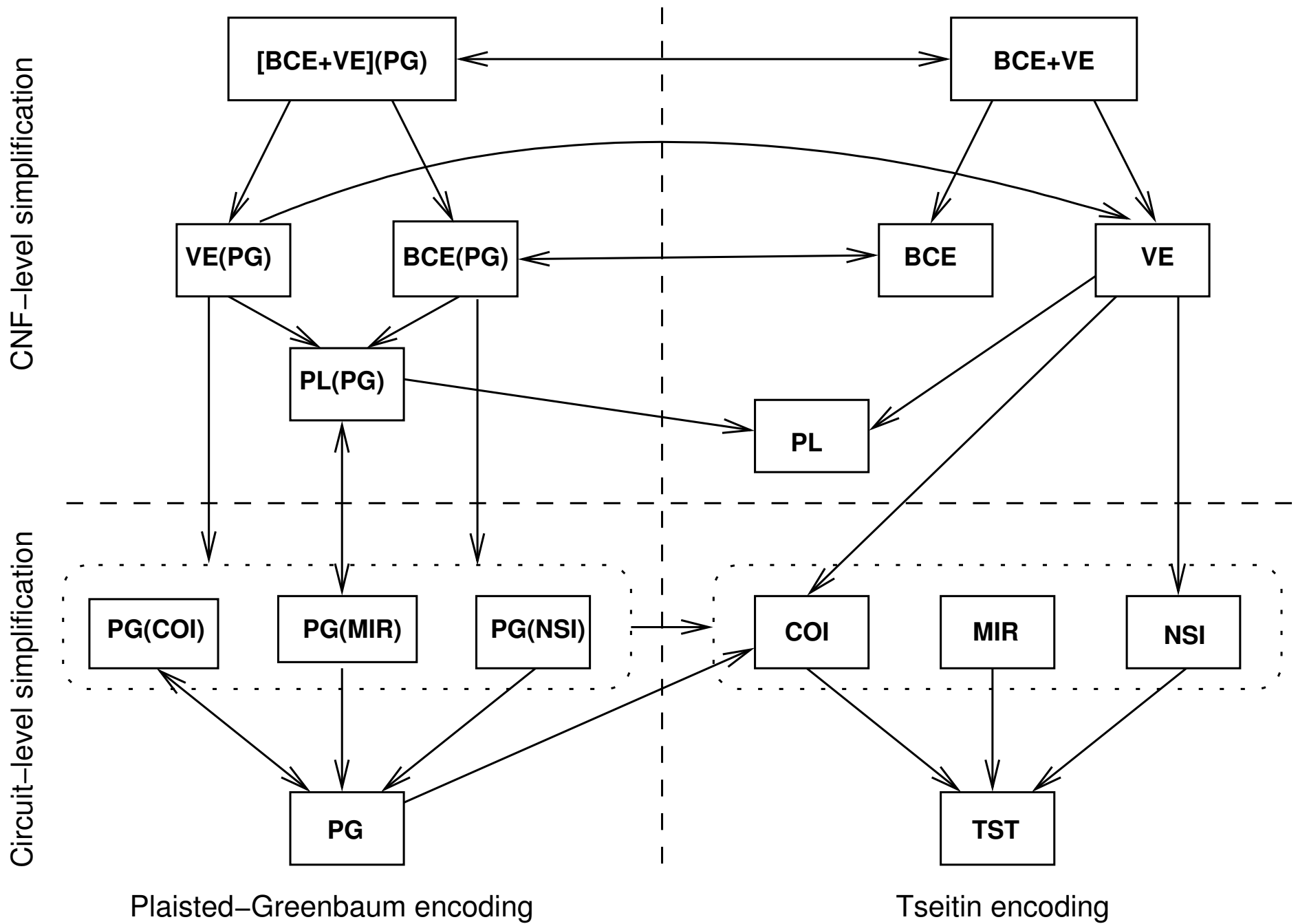
**PG** Plaisted-Greenbaum polarity based encoding

**TST** standard Tseitin encoding

---

**VE** Variable-Elimination as in DP / Quantor / SATeLite

**BCE** Blocked-Clause-Elimination



**[BCE+VE](PG)**

**BCE+VE**

**VE(PG)**

**BCE(PG)**

**BCE**

**VE**

**PL(PG)**

**PL**

**PG(COI)**

**PG(MIR)**

**PG(NSI)**

**COI**

**MIR**

**NSI**

**PG**

**TST**

Plaisted-Greenbaum encoding

Tseitin encoding

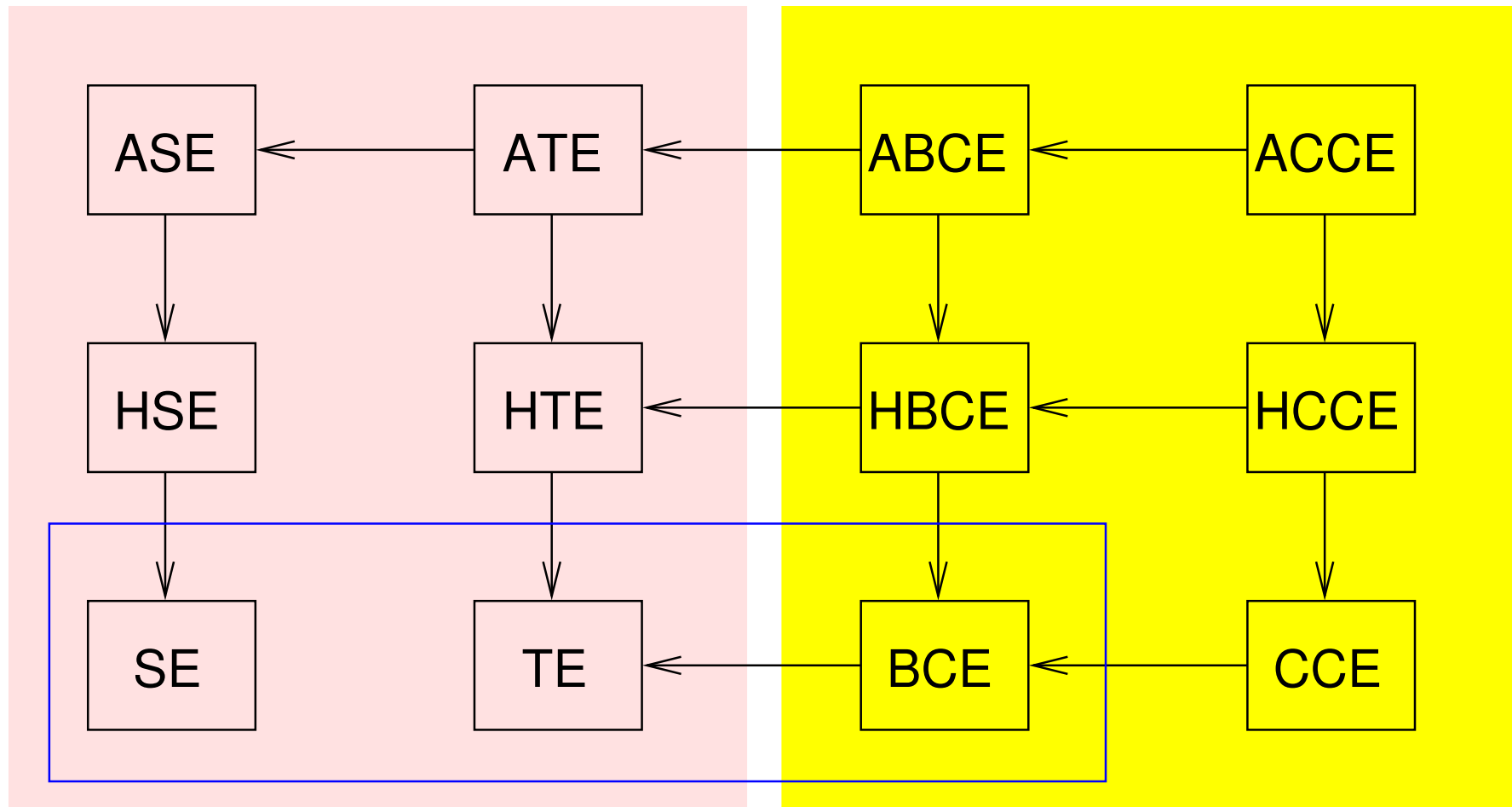
	encoding			b			be			beb			bebe			e		
	T	V	C	T	V	C	T	V	C	T	V	C	T	V	C	T	V	C
SU	0	46	256	2303	29	178	1042	11145		1188	11145		569	11144		2064	11153	
AT	12	9	27	116	7	18	1735	1	8	1835	1	6	34	1	6	244	1	9
AP	10	9	20	94	7	18	1900	1	6	36	1	6	34	1	6	1912	1	6
AM	190	1	8	42	1	7	178	1	7	675	1	7	68	1	7	48	1	8
AN	9	3	10	50	3	10	1855	1	6	36	1	6	34	1	6	1859	1	6
HT	147	121	347	1648	117	277	2641	18118		567	18118		594	18116		3240	23140	
HP	130	121	286	1398	117	277	2630	18118		567	18118		595	18116		2835	19119	
HM	6961	16	91	473	16	84	621	12	78	374	12	77	403	12	76	553	15	90
HN	134	34	124	573	34	122	1185	17102		504	17101		525	17100		1246	17103	
BT	577	442	1253	5799	420	1119	7023	57321		1410	56310		1505	52294		8076	64363	
BP	542	442	1153	5461	420	1119	7041	57321		1413	56310		1506	52294		7642	57322	
BM	10024	59	311	1252	58	303	1351	53287		1135	53286		1211	52280		1435	55303	
BN	13148	196	643	2902	193	635	4845	108508		2444	107504		2250	105500		5076	114518	

S = Sat competition  
 A = AIG competition  
 H = HW model checking competition  
 B = bit-vector SMT competition

T = plain Tseitin encoding  
 P = Plaisted Greenbaum  
 M = MiniCirc encoding  
 N = NiceDAGs



H = hidden, A = asymmetric,  
 SE = subsumption elimination, T = tautology elimination  
 BC = blocked clause elimination, CC = covered clause elimination



logically equivalent

satisfiability equivalent

**Definition** [self-subsuming resolution]

[EénBiere SAT'05]

if  $A \vee x$  and  $B \vee \bar{x}$  are two clauses in a CNF and  $\frac{A \vee x \quad \bar{x} \vee B}{A}$  a valid resolution, e.g.  $B \subseteq A$ , then replace the clause  $A \vee x$  by  $A$ , in essence **removing**  $x$  from  $A \vee x$

example: if both  $a \vee b \vee x$  and  $b \vee \bar{x}$  are in a CNF remove  $x$  from first clause

**Definition** [asymmetric literal addition]

[HeuleJärvisaloBiere LPAR'10]

if  $A$  and  $B \vee \bar{x}$  are two clauses in a CNF and  $\frac{A \vee x \quad \bar{x} \vee B}{A}$  a valid resolution, e.g.  $B \subseteq A$ , then replace the clause  $A$  by  $A \vee x$ , in essence **adding**  $x$  to  $A$

example: if both  $a \vee b$  and  $b \vee \bar{x}$  are in a CNF add  $x$  to first clause

**Definition** [asymmetric tautology / blocked clause] [HeuleJärvisaloBiere LPAR'10]

apply asymmetric literal addition to a clause w.r.t. to fixed CNF as long as possible, if the result is a tautological / blocked then remove clause (otherwise keep original)

**Definition** hidden = only use binary side clauses  $\bar{x} \vee B$

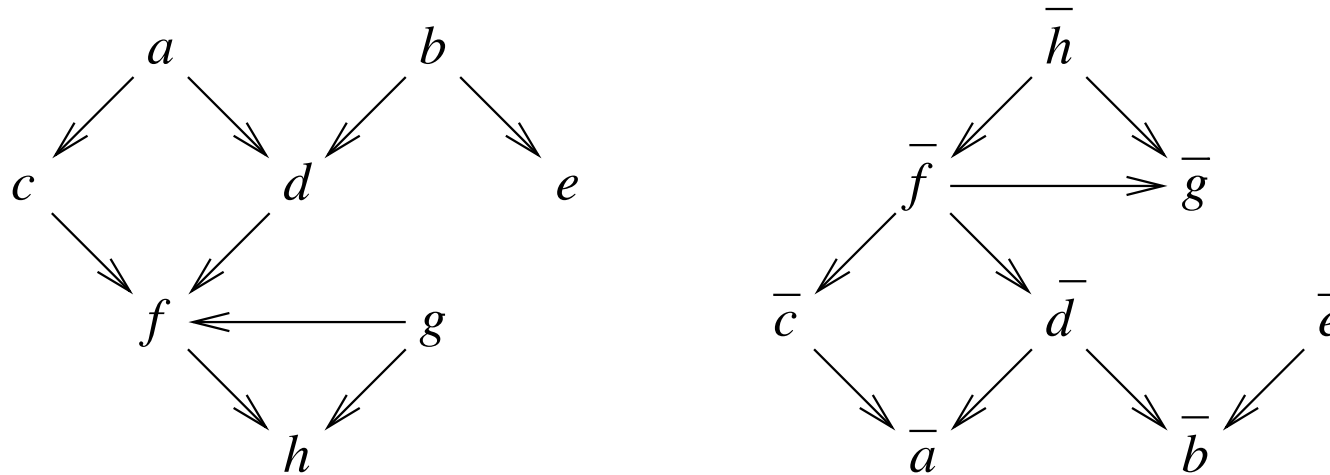
**Fact** AHTE can be simulated by asymmetric branching / distillation and BCP

**Fact** HTE can be implemented much faster by iterating over all literals instead of iterating over all clauses (partially implemented in Lingeling)

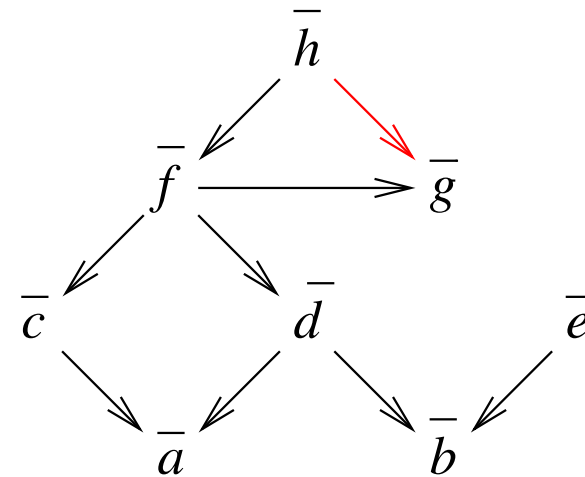
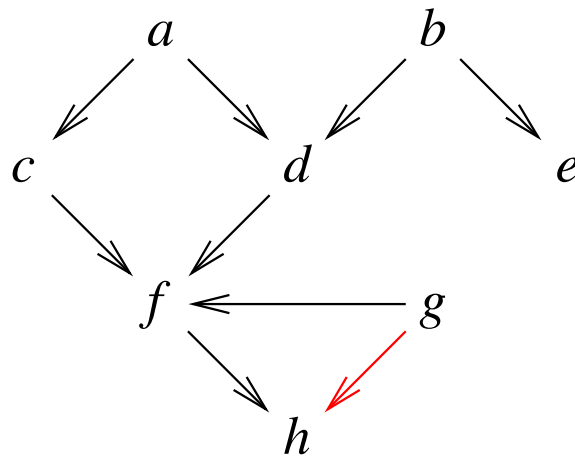
**Fact** HTE after equivalence reasoning and failed literal probing until completion on binary clauses only is confluent and BCP preserving

see our long and short LPAR'10 papers for more details

- SAT solvers applied to huge formulas
  - million of variables
  - fastest solvers use preprocessing/inprocessing
  - *need cheap and effective inprocessing techniques for millions of variables*
- this talk:
  - **unhiding** redundancy in large formulas
  - almost linear randomized algorithm
  - using the binary implication graph
  - fast enough to be applied to learned clauses
- see our SAT'11 paper for more details



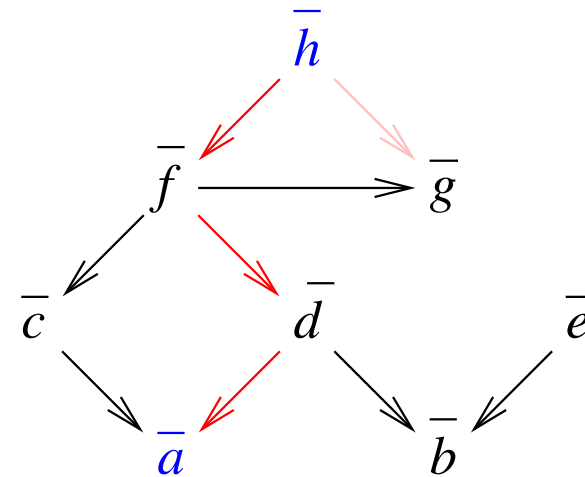
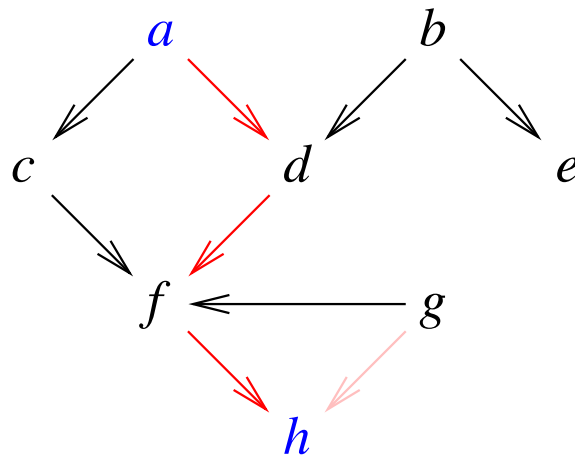
$$\begin{aligned}
 & (\bar{a} \vee c) \wedge (\bar{a} \vee d) \wedge (\bar{b} \vee d) \wedge (\bar{b} \vee e) \wedge \\
 & (\bar{c} \vee f) \wedge (\bar{d} \vee f) \wedge (\bar{g} \vee f) \wedge (\bar{f} \vee h) \wedge \\
 & (\bar{g} \vee h) \wedge \underbrace{(\bar{a} \vee \bar{e} \vee h) \wedge (\bar{b} \vee \bar{c} \vee h) \wedge (a \vee b \vee c \vee d \vee e \vee f \vee g \vee h)}_{\text{non binary clauses}}
 \end{aligned}$$



$$\begin{aligned}
 & (\bar{a} \vee c) \wedge (\bar{a} \vee d) \wedge (\bar{b} \vee d) \wedge (\bar{b} \vee e) \wedge \\
 & (\bar{c} \vee f) \wedge (\bar{d} \vee f) \wedge (\bar{g} \vee f) \wedge (\bar{f} \vee h) \wedge \\
 & \cancel{(\bar{g} \vee h)} \wedge (\bar{a} \vee \bar{e} \vee h) \wedge (\bar{b} \vee \bar{c} \vee h) \wedge (a \vee b \vee c \vee d \vee e \vee f \vee g \vee h)
 \end{aligned}$$

**TRD**

$$g \rightarrow f \rightarrow h$$



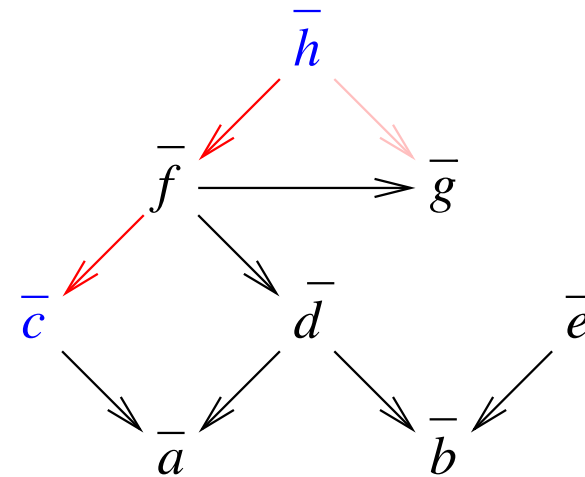
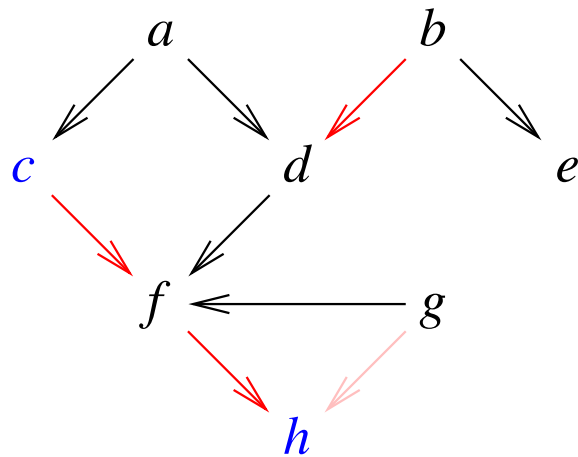
$$(\bar{a} \vee c) \wedge (\bar{a} \vee d) \wedge (\bar{b} \vee d) \wedge (\bar{b} \vee e) \wedge$$

$$(\bar{c} \vee f) \wedge (\bar{d} \vee f) \wedge (\bar{g} \vee f) \wedge (\bar{f} \vee h) \wedge$$

$$(\bar{a} \vee \bar{e} \vee h) \wedge (\bar{b} \vee \bar{c} \vee h) \wedge (a \vee b \vee c \vee d \vee e \vee f \vee g \vee h)$$

**HTE**

$$a \rightarrow d \rightarrow f \rightarrow h$$



$$(\bar{a} \vee c) \wedge (\bar{a} \vee d) \wedge (\bar{b} \vee d) \wedge (\bar{b} \vee e) \wedge$$

$$(\bar{c} \vee f) \wedge (\bar{d} \vee f) \wedge (\bar{g} \vee f) \wedge (\bar{f} \vee h) \wedge$$

$$~~(\bar{b} \vee \bar{c} \vee h)~~ \wedge (a \vee b \vee c \vee d \vee e \vee f \vee g \vee h)$$

**HTE**

$$c \rightarrow f \rightarrow h$$



$$\frac{C \vee l \quad D \vee \bar{l}}{D} \quad C \subseteq D$$

$$\frac{a \vee b \vee l \quad a \vee b \vee c \vee \bar{l}}{a \vee b \vee c}$$

resolvent  $D$  subsumes second antecedent  $D \vee \bar{l}$

assume given CNF contains both antecedents

$$\dots (a \vee b \vee l) (a \vee b \vee c \vee \bar{l}) \dots$$

if  $D$  is added to CNF then  $D \vee \bar{l}$  can be removed

$$\Downarrow$$

which in essence *removes*  $\bar{l}$  from  $D \vee \bar{l}$

$$\dots (a \vee b \vee l) (a \vee b \vee c) \dots$$

used in SATeLite preprocessor

now common in many SAT solvers

hidden literal addition (HLA) uses SSR in reverse order

$$\frac{C \vee l \quad D \vee \bar{l}}{D} \quad C \subseteq D$$

$$\frac{a \vee b \vee l \quad a \vee b \vee c \vee \bar{l}}{a \vee b \vee c}$$

assume given CNF contains resolvent and first antecedent

$$\dots (a \vee b \vee l)(a \vee b \vee c) \dots$$

we can replace  $D$  by  $D \vee \bar{l}$

$$\dots (a \vee b \vee l)(a \vee b \vee c \vee \bar{l}) \dots$$

which in essence *adds*  $\bar{l}$  to  $D$ , repeat HLA until fix-point

keep remaining non-tautological clauses *after removing added literals again*

HTE = assume  $C \vee l$  is a binary clauses

more general versions in the paper

**remove clauses with a literal implied by negation of another literal in the clause**

HTE confluent and BCP preserving

modulo equivalent variable renaming

better explained on binary implication graph

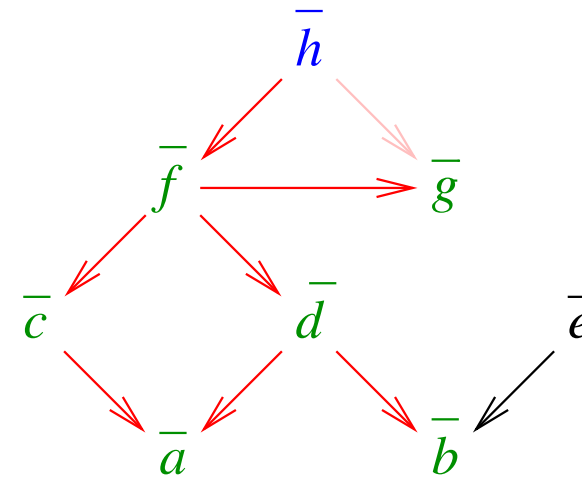
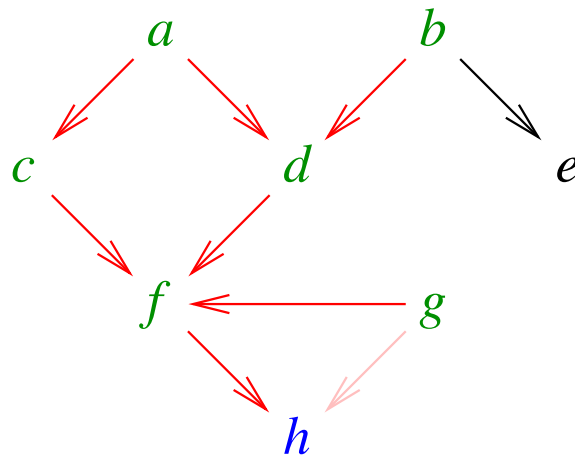
**remove literal from a clause which implies another literal in the clause**

$$\dots(\bar{a} \vee b)(\bar{b} \vee c)(a \vee c \vee d)\dots \Rightarrow \dots(\bar{a} \vee b)(\bar{b} \vee c)(c \vee d)\dots$$

related work before all uses BCP:

- asymmetric branching implemented in MiniSAT but switched off by default
- **distillation** [JinSomenzi'05][HanSomenzi DAC'07]
- vivification [PietteHamadiSais ECAI'08]
- caching technique in CryptoMiniSAT

HTE/HLE only uses the binary implication graph!



$$(\bar{a} \vee c) \wedge (\bar{a} \vee d) \wedge (\bar{b} \vee d) \wedge (\bar{b} \vee e) \wedge$$

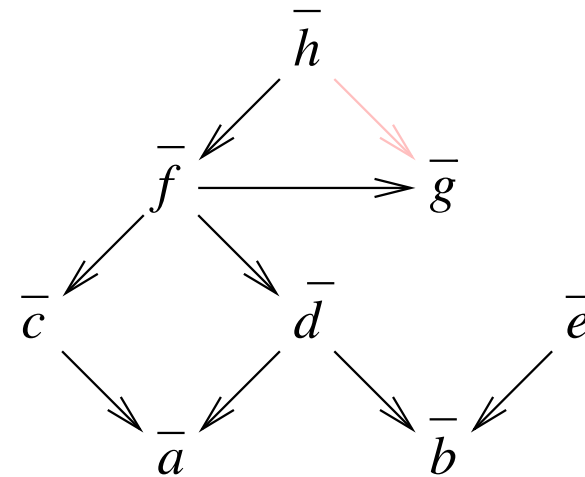
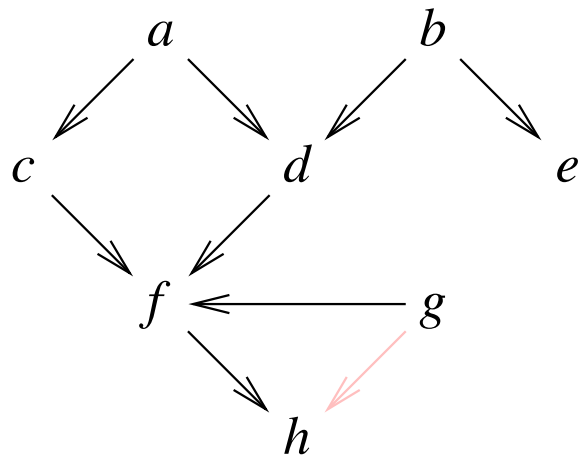
$$(\bar{c} \vee f) \wedge (\bar{d} \vee f) \wedge (\bar{g} \vee f) \wedge (\bar{f} \vee h) \wedge$$

$$(\cancel{a} \vee \cancel{b} \vee \cancel{c} \vee \cancel{d} \vee e \vee \cancel{f} \vee \cancel{g} \vee h)$$

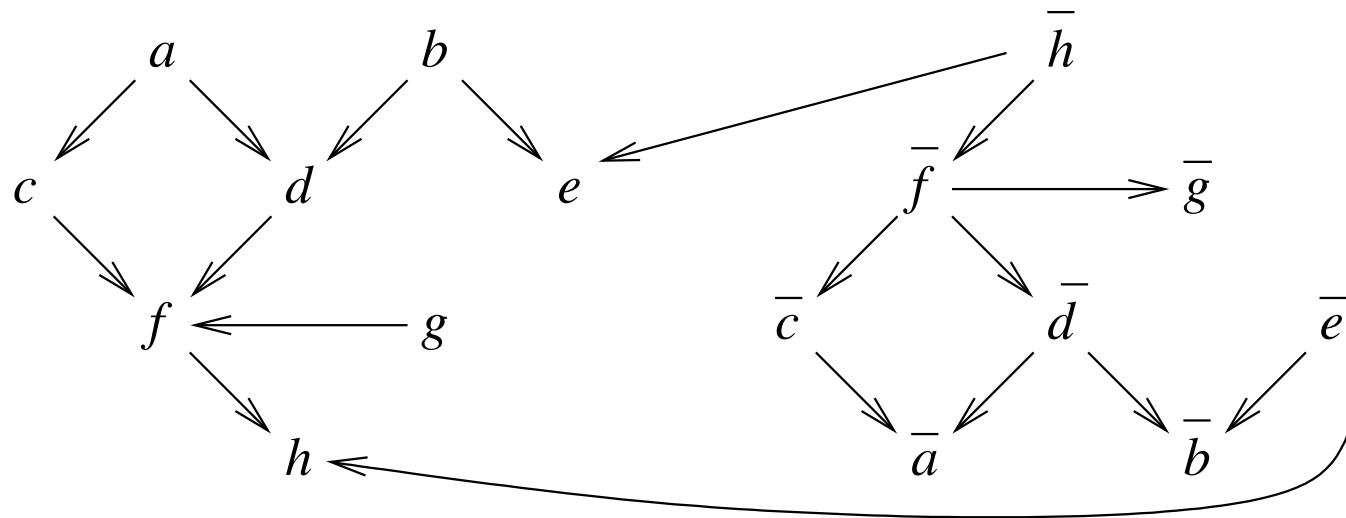
**HLE**

all but  $e$  imply  $h$

also  $b$  implies  $e$



$$\begin{aligned}
 & (\bar{a} \vee c) \wedge (\bar{a} \vee d) \wedge (\bar{b} \vee d) \wedge (\bar{b} \vee e) \wedge \\
 & (\bar{c} \vee f) \wedge (\bar{d} \vee f) \wedge (\bar{g} \vee f) \wedge (\bar{f} \vee h) \wedge \\
 & ( \qquad \qquad \qquad e \vee \qquad \qquad \qquad h )
 \end{aligned}$$



$$\begin{aligned}
 &(\bar{a} \vee c) \wedge (\bar{a} \vee d) \wedge (\bar{b} \vee d) \wedge (\bar{b} \vee e) \wedge \\
 &(\bar{c} \vee f) \wedge (\bar{d} \vee f) \wedge (\bar{g} \vee f) \wedge (\bar{f} \vee h) \wedge \\
 &(e \vee h)
 \end{aligned}$$

actually quite old technique

... [Freeman PhdThesis'95] [LeBerre'01] ...

*assume* literal  $l$ , BCP, if conflict, add unit  $\bar{l}$

rather costly to run until completion

conjecture: at least quadratic

one BCP is linear and also in practice can be quite expensive

need to do it for all variables and restart if new binary clause generated

useful in practice: lift common implied literals for assumption  $l$  and assumption  $\bar{l}$

**even on BIG (FL2) conjectured to be quadratic**

[VanGelder'05]

...  $(\bar{a} \vee b)(\bar{b} \vee c)(\bar{c} \vee d)(\bar{d} \vee \bar{a})$  ...  $\Rightarrow$  add unit clause  $\bar{a}$

subsumed by running one HLA until completion

decompose BIG into strongly connect components (SCCs)

if there is an  $l$  with  $l$  and  $\bar{l}$  in the same component  $\Rightarrow$  *unsatisfiable*

otherwise replace all literals by a “representative”

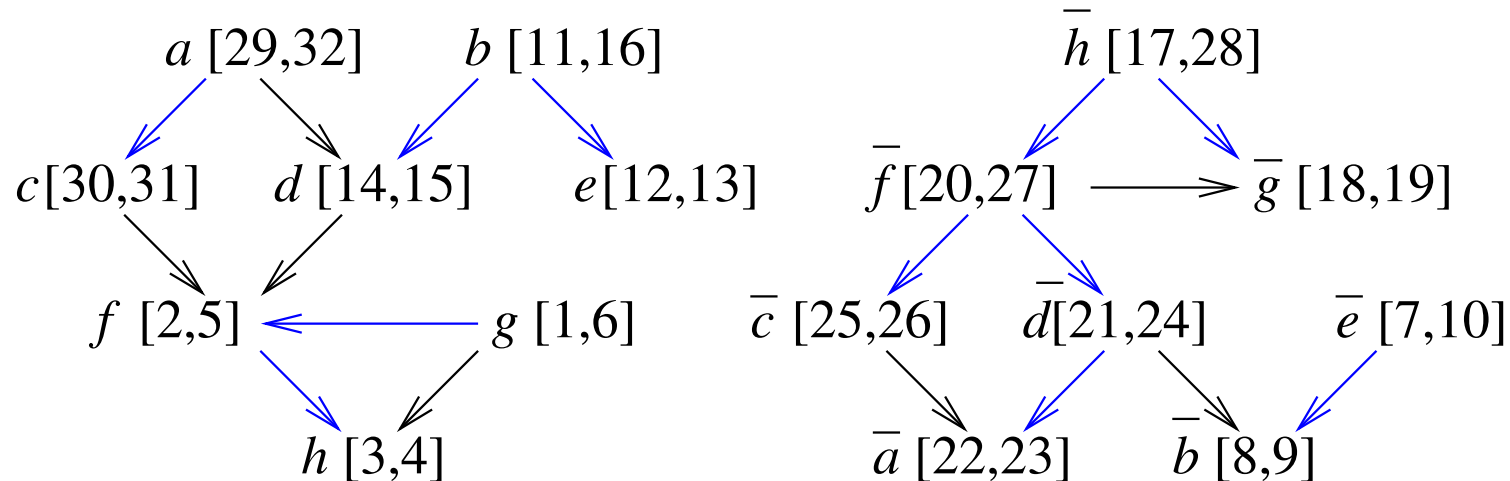
**linear algorithm** can be applied routinely during garbage collection

but as with failed literal preprocessing may generate new binary clauses

$$\dots(\bar{a} \vee b)(\bar{b} \vee c)(\bar{c} \vee a)(a \vee b \vee c \vee d)\dots \Rightarrow \dots(a \vee d)\dots$$



DFS tree with discovered and finished times:  $[dsc(l), fin(l)]$



tree edges

parenthesis theorem:  $l$  ancestor in DFS tree of  $k$  iff  $[dsc(k), fin(k)] \subseteq [dsc(l), fin(l)]$   
 well known

ancestor relationship gives necessary conditions for (transitive) implication:

if  $[dsc(k), fin(k)] \subseteq [dsc(l), fin(l)]$  then  $l \rightarrow k$

if  $[dsc(\bar{l}), fin(\bar{l})] \subseteq [dsc(\bar{k}), fin(\bar{k})]$  then  $l \rightarrow k$

- time stamping in previous example does not cover  $b \rightarrow h$   
 $[11, 16] = [\text{dsc}(b), \text{fin}(b)] \not\subseteq [\text{dsc}(h), \text{fin}(h)] = [3, 4]$   
 $[17, 28] = [\text{dsc}(\bar{h}), \text{fin}(\bar{h})] \not\subseteq [\text{dsc}(\bar{b}), \text{fin}(\bar{b})] = [8, 9]$
- in example still both HTE “unhidden”, HLE works too (since  $b \rightarrow e$ )
- “coverage” heavily depends on DFS order
- as solution we propose multiple **randomized DFS** rounds/phases
- so we approximate a quadratic problem (reachability) randomly by a linear algorithm
- if BIG is a tree *one* time stamping covers everything

*Unhiding* (formula  $F$ )

```

1  stamp := 0
2  foreach literal  $l$  in  $BIG(F)$  do
3     $dsc(l) := 0$ ;  $fin(l) := 0$ 
4     $prt(l) := l$ ;  $root(l) := l$ 
5  foreach  $r \in RTS(F)$  do
6    stamp := Stamp( $r, stamp$ )
7  foreach literal  $l$  in  $BIG(F)$  do
8    if  $dsc(l) = 0$  then
9      stamp := Stamp( $l, stamp$ )
10 return Simplify( $F$ )

```

*Stamp* (literal  $l$ , integer *stamp*)

```

1  stamp := stamp + 1
2   $dsc(l) := stamp$ 
3  foreach  $(\bar{l} \vee l') \in F_2$  do
4    if  $dsc(l') = 0$  then
5       $prt(l') := l$ 
6       $root(l') := root(l)$ 
7      stamp := Stamp( $l', stamp$ )
8  stamp := stamp + 1
9   $fin(l) := stamp$ 
10 return stamp

```

*Simplify* (formula  $F$ )

```

1  foreach  $C \in F$ 
2     $F := F \setminus \{C\}$ 
3    if  $UHTE(C)$  then continue
4     $F := F \cup \{UHLE(C)\}$ 
5  return  $F$ 

```

*UHTE* (clause  $C$ )

```

1    $l_{\text{pos}} :=$  first element in  $S^+(C)$ 
2    $l_{\text{neg}} :=$  first element in  $S^-(C)$ 
3   while true
4     if  $\text{dsc}(l_{\text{neg}}) > \text{dsc}(l_{\text{pos}})$  then
5       if  $l_{\text{pos}}$  is last element in  $S^+(C)$  then return false
6        $l_{\text{pos}} :=$  next element in  $S^+(C)$ 
7     else if  $\text{fin}(l_{\text{neg}}) < \text{fin}(l_{\text{pos}})$  or ( $|C| = 2$  and ( $l_{\text{pos}} = \bar{l}_{\text{neg}}$  or  $\text{prt}(l_{\text{pos}}) = l_{\text{neg}}$ )) then
8       if  $l_{\text{neg}}$  is last element in  $S^-(C)$  then return false
9        $l_{\text{neg}} :=$  next element in  $S^-(C)$ 
10  else return true

```

$S^+(C)$  sequence of literals in  $C$  ordered by  $\text{dsc}()$

$S^-(C)$  sequence of negations of literals in  $C$  ordered by  $\text{dsc}()$

$$O(|C|\log|C|)$$

*UHLE* (clause  $C$ )

```
1  finished := finish time of first element in  $S_{\text{rev}}^+(C)$ 
2  foreach  $l \in S_{\text{rev}}^+(C)$  starting at second element
3      if  $\text{fin}(l) > \textit{finished}$  then  $C := C \setminus \{l\}$ 
4      else finished :=  $\text{fin}(l)$ 
5  finished := finish time of first element in  $S^-(C)$ 
6  foreach  $\bar{l} \in S^-(C)$  starting at second element
7      if  $\text{fin}(\bar{l}) < \textit{finished}$  then  $C := C \setminus \{l\}$ 
8      else finished :=  $\text{fin}(\bar{l})$ 
9  return  $C$ 
```

$S_{\text{rev}}^+(C)$  reverse of  $S^+(C)$

$O(|C|\log|C|)$

```

Stamp (literal  $l$ , integer  $stamp$ )
1 BSC    $stamp := stamp + 1$ 
2 BSC    $dsc(l) := stamp; obs(l) := stamp$ 
3 ELS    $flag := true$  //  $l$  represents a SCC
4 ELS    $S.push(l)$  // push  $l$  on SCC stack
5 BSC   for each  $(\bar{l} \vee l') \in F_2$ 
6 TRD   if  $dsc(l) < obs(l')$  then  $F := F \setminus \{(\bar{l} \vee l')\}$ ; continue
7 FLE   if  $dsc(\text{root}(l)) \leq obs(\bar{l}')$  then
8 FLE    $l_{\text{failed}} := l$ 
9 FLE   while  $dsc(l_{\text{failed}}) > obs(\bar{l}')$  do  $l_{\text{failed}} := \text{prt}(l_{\text{failed}})$ 
10 FLE   $F := F \cup \{(\bar{l}_{\text{failed}})\}$ 
11 FLE  if  $dsc(\bar{l}')$   $\neq 0$  and  $\text{fin}(\bar{l}')$   $= 0$  then continue
12 BSC  if  $dsc(l')$   $= 0$  then
13 BSC   $\text{prt}(l') := l$ 
14 BSC   $\text{root}(l') := \text{root}(l)$ 
15 BSC   $stamp := \text{Stamp}(l', stamp)$ 
16 ELS  if  $\text{fin}(l')$   $= 0$  and  $dsc(l')$   $< dsc(l)$  then
17 ELS   $dsc(l) := dsc(l')$ ;  $flag := false$  //  $l$  is equivalent to  $l'$ 
18 OBS   $obs(l') := stamp$  // set last observed time attribute
19 ELS  if  $flag = true$  then // if  $l$  represents a SCC
20 BSC   $stamp := stamp + 1$ 
21 ELS  do
22 ELS   $l' := S.pop()$  // get equivalent literal
23 ELS   $dsc(l') := dsc(l)$  // assign equal discovered time
24 BSC   $\text{fin}(l') := stamp$  // assign equal finished time
25 ELS  while  $l' \neq l$ 
26 BSC  return  $stamp$ 

```

- implemented as one inprocessing phase in our SAT solver Lingeling  
beside variable elimination, distillation, blocked clause elimination, probing, ...
- bursts of randomized DFS rounds and sweeping over the whole formula
- fast enough to be applicable to large learned clauses as well  
unhiding is particularly effective for learned clauses
- beside UHTE and UHLE we also have added hyper binary resolution UHBR  
not useful in practice

configuration	sol	sat	uns	unhd	simp	elim
adv.stamp (no uhbr)	188	78	110	7.1%	33.0%	16.1%
adv.stamp (w/uhbr)	184	75	109	7.6%	32.8%	15.8%
basic stamp (no uhbr)	183	73	110	6.8%	32.3%	15.8%
basic stamp (w/uhbr)	183	73	110	7.4%	32.8%	15.8%
no un hiding	180	74	106	0.0%	28.6%	17.6%

configuration	hte	stamp	redundant	hle	redundant	units	stamp
adv.stamp (no uhbr)	22	64%	59%	291	77.6%	935	57%
adv.stamp (w/uhbr)	26	67%	70%	278	77.9%	941	58%
basic stamp (no uhbr)	6	0%	52%	296	78.0%	273	0%
basic stamp (w/uhbr)	7	0%	66%	288	76.7%	308	0%
no un hiding	0	0%	0%	0	0.0%	0	0%

similar results for crafted and SAT'10 Race instances



- preprocessing is important for SAT solvers
- hard kernels do occur in practice
- inprocessing provides additional benefits
- new class of clause elimination procedures
- even quadratic algorithms are most of the time too expensive