# On the representation and analysis of distributed search in theorem proving

MARIA PAOLA BONACINA

DEPT. OF COMPUTER SCIENCE

THE UNIVERSITY OF IOWA

# Research program

Area: Automated Reasoning

Emphasis: Control of Deduction

Directions:

* Combination of forward and backward
  reasoning, e.g.,

  Target - oriented equational reasoning

  Lemmatization in semantic strategies

* Distributed automated deduction, e.g.,

  Clause- Diffusion methodology

  Modified Clause- Diffusion

  AGO - criteria

  Combination of distributed search and multi-search

  Systems built: Aquarius, Peers, Peers-mcd

* Strategy analysis, e.g.,

  Search space reduction by contraction

  Distributed search for contraction-based strategies

Motivation:

- Applications
- Impact on other areas in C.S.
- Basic investigation

# Theorem proving and parallelism
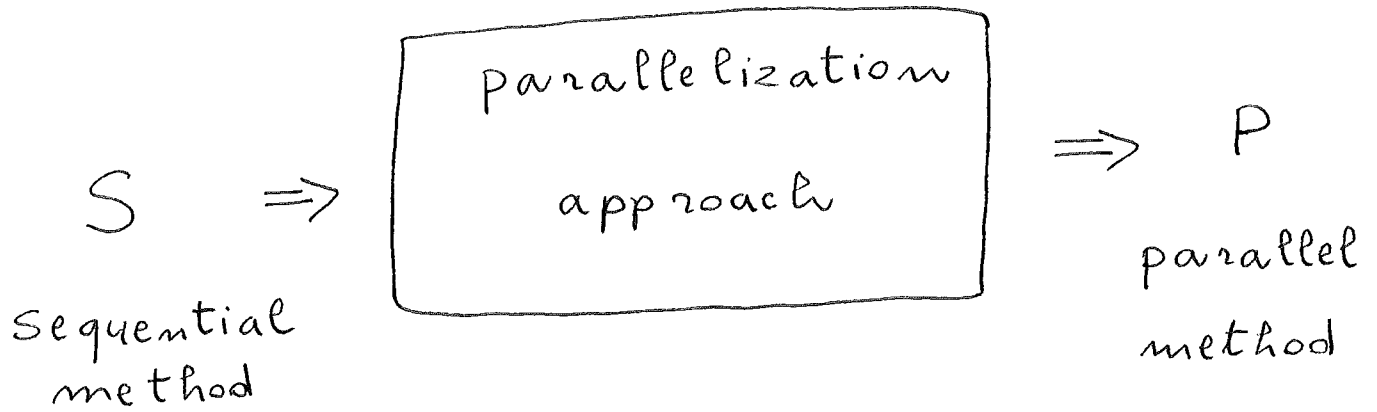
- More power:

  faster proofs

  more proofs


- Search plan design:

  investigation of new forms

  of control of deduction

## Many approaches:

$$S \Rightarrow \boxed{\begin{array}{c} \text{parallelization} \\ \text{approach} \end{array}} \Rightarrow P$$

Sequential method

parallel method

## Evaluation:

- performance evaluation

- analysis ?

# Outline

Contraction - based T.P. strategies
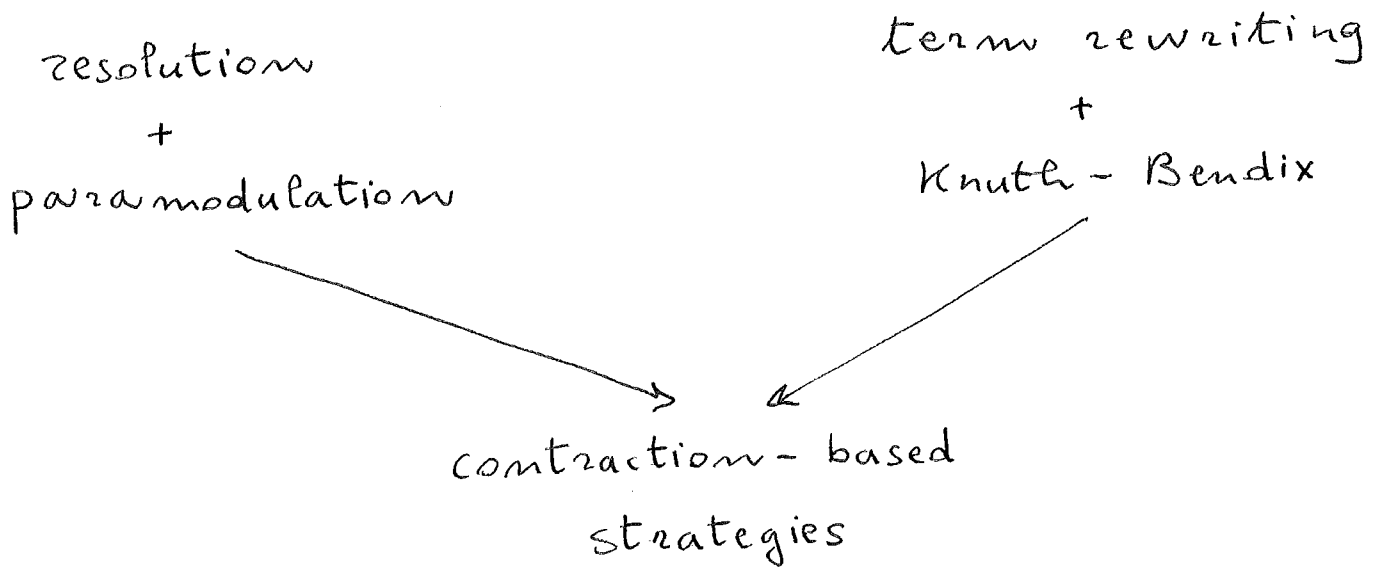
Distributed search

Representation

Analysis

Comparison of distributed strategy and sequential base

Discussion

What are contraction - based strategies and why are they important ?

# Contraction- based strategies

History:

resolution
+
paramodulation

term rewriting
+
Knuth - Bendix

contraction- based
strategies

Forward reasoning:
   generate and keep clauses

Ordering - based:
   well founded $<$ on Terms and clauses

Good for equality reasoning

# Contraction-based strategies

$$\mathcal{C} = \langle I, \Sigma \rangle$$

I: inference rules

Expansion     (e.g., resolution)

Contraction     (e.g., simplification)

$\Sigma$: search plan

$$\Sigma = \langle \varsigma, \xi, \omega \rangle$$

$\varsigma$: select rule     $[\varsigma: States^* \longrightarrow I]$

$\xi$: select premises     $[\xi: States^* \longrightarrow \mathcal{L}]$

$\omega$: detect success     $[\omega: States \longrightarrow Boolean]$

$I$ refutationally complete $\left.\begin{array}{c} \\ \\ \end{array}\right\}$ $\mathcal{C}$ complete

$\Sigma$ fair

$$S_0 \vdash S_1 \vdash \ldots S_i \vdash \ldots$$

# Contraction - based strategies

**Forward contraction**

normalize every new clauses w.r.t. existing ones

**Backward contraction**

normalize every existing clause w.r.t. new insertions $\Rightarrow$ inter - reduction

**Eager contraction**

$\mathcal{E}$ does not select expansion until contraction exhausted

$S_0 \vdash \ldots S_i \vdash \ldots \qquad \forall i \qquad \forall \varphi \in S_i$

if $\exists f \quad \exists \bar{x} \in S_i \quad f$ applied to $(\bar{x}, \varphi)$ deletes $\varphi$

then $\exists \ell \geq i \quad S_\ell \vdash S_{\ell+1}$ deletes $\varphi$ and

$\forall j \quad i \leq j \leq \ell$ no expansion unless succeeds

sooner

# Some results of CBS

- Moufang identities in rings
  Anantharaman, Hsiang     SBR 2    1990

- Axiomatization of Lukasiewicz many-valued logic
  Anantharaman, Bonacina    SBR 3    1989-91

- Single axioms for groups
  W. McCune     OTTER    1993

- Robbins algebras are Boolean
  W. McCune     EQP    1996

- Verification of cryptographic protocols
  C. Weidenbach     SPASS    1999

What is distributed search and why do we use it ?

# Parallelism at the term level

Parallelize the single inference
(e.g., parallel rewriting)

Motivation: speed-up frequent operations

Good for: concurrent rewriting

Not for CBS:

- New equations generated dynamically

$$\Downarrow$$

  special pre-processing too onerous

- Many terms, equations, steps:
  too fine-grained

## Parallelism at the clause level

Parallel inferences within one search
(e.g., parallel resolution steps,
OR - parallelism)

Motivation: speed-up given search

Good for: expansion - oriented T.P.
(e.g., hyperresolution with
no contraction)

Not for CBS:

backward contraction causes conflicts

⇓

do it sequentially : bottleneck

# Parallelism at the search level

Parallel derivations:
  deductive processes search in parallel
  the space of the problem


Heterogeneous systems:
  different inference systems
  motivation: combine forward / backward
                        reasoning


Multi-search:
  different search plans
  motivation: search in different order


Distributed search:
  subdivide search space
  motivation: divide work


All need communication

# Distributed strategy

$$\ell = \langle I, M, \Sigma \rangle$$

$I$: inference rules
(expansion rules, contraction rules)

$M$: communication operators
(send, receive ...)

$\Sigma$: search plan
    sequential:      select rule
                           select premises

    distributed:      also subdivision
                               communication

Distributed-search plan

$$\Sigma = \langle \zeta, \xi, \alpha, \omega \rangle$$

$\zeta$: select rule / operator

$\xi$: select premises

$\alpha$: subdivision function

$\omega$: detects termination

# Subdivision function α

Subdivide inferences among $P_0 \, P_1 \, \ldots \, P_{m-1}$

Search space infinite unknown $\Rightarrow$
dynamic subdivision : at each stage $S_i$
of derivation subdivide inferences in $S_i$

$$\alpha \left( S_0 \ldots S_i , \, m, \, \kappa, \, f, \, \bar{x} \right) = \text{true} \, / \, \text{false}$$

means $P_\kappa$ is allowed / forbidden
to apply $f$ to $\bar{x}$

Two requirements :

$\alpha$ is total on generated clauses
( partial function in general)

$\alpha$ monotonic : $\qquad \left( \text{w.r.t.} \; i \right)$

$$
\begin{array}{ccc|ccc}
\bot & \bot & \bot & \text{false} & \text{false} & \text{false} & \ldots \\
\bot & \bot & \bot & \text{true} & \text{true} & \text{true} & \ldots \\
\bot & \bot & \bot & \bot & \bot & \bot & \ldots \\
\end{array}
$$

# Parallelization by subdivision

$$\ell = \langle I, \mathcal{E} \rangle \qquad \mathcal{E} = \langle \zeta, \xi, \omega \rangle$$

$$\ell' = \langle I, M, \mathcal{E}' \rangle \qquad \mathcal{E}' = \langle \zeta', \xi', \alpha, \omega \rangle$$

$\zeta'$ and $\xi'$ select inferences like $\zeta$ and $\xi$ so that the difference is made by $\alpha$: forbidden steps $\Rightarrow$ different selections and by the presence of communication.

$P_0, P_1 \ldots P_{m-1}$ :  different derivations:

$$S = S_0^0 \vdash S_1^0 \vdash \ldots S_i^0 \vdash \ldots$$
$$\vdots$$
$$S = S_0^k \vdash S_1^k \vdash \ldots S_i^k \vdash \ldots$$
$$\vdots$$
$$S = S_0^{m-1} \vdash S_1^{m-1} \vdash \ldots S_i^{m-1} \vdash \ldots$$

$\left. \right\}$ distributed derivation

# Fairness of distributed derivations

Refutational completeness of $I$ +
Fairness of $\Sigma$ = Completeness of $C$

$\forall \bar{x}$    persistent    non-redundant
$\forall f$    expansion    rule
$\exists P_k$    such that

(1)   $P_k$ has $\bar{x}$    (fairness of communication)

(2)   $P_k$ is allowed to apply $f$ to $\bar{x}$
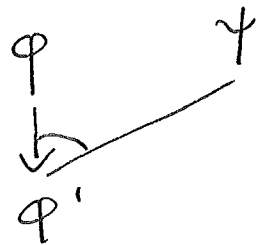     at some stage    (fairness of subdivision)

and (3) all local derivations are fair
                         (local fairness)

Theorem:      (1) + (2) + (3)     $\Rightarrow$

the distributed derivation is fair

How can we guarantee that a parallelization by subdivision of a contraction - based strategy is also contraction - based ?

# Eager contraction in distributed derivation

$$\varphi \searrow \psi \quad\quad \text{in } \bigcup_h S_i^h$$
$$\varphi'$$

What if $\varphi \in S_i^n$ and $\psi \in S_i^j$ ?

What if the step is forbidden?

## Propagation of clauses up to redundancy:

$\varphi$ persistent non-redundant

$\varphi \in \bigcup_h S_i^h$ (i first stage)

then $\forall p_h \; \exists j \; \varphi \in S_j^h$ (delay: $j - i \geqslant 0$)

/* also sufficient for fairness of communication */

# Eager contraction in distributed derivation

Distributed global contraction:
$$\forall P_\kappa \quad \forall i \quad \forall \varphi \in S_i^\kappa$$
if $\exists f \; \exists \bar{x} \in \bigcup_\kappa S_i^\ell$ $f$ applied to $\bar{x}$ deletes $\varphi$
then $\exists \ell \geqslant i \; P_\kappa$ deletes $\varphi$ at stage $\ell$
unless it halts sooner.

Global eager contraction:
no expansion      in between     $(\forall j \; i \leqslant j \leqslant \ell)$
no communication

## Lemmas:

Local eager
contraction

Propagation of clauses
up to redundancy

$\Big\} \Rightarrow$ Distributed
global
contraction

Local eager
contraction

Immediate propagation
of clauses up to
redundancy

$\Big\} \Rightarrow$ Global
eager
contraction

# Contraction - based strategies

**Sequential** :     contraction rules

eager contraction

**Distributed** :     contraction rules

distributed global contraction

$\mathcal{C}$ :   contraction - based

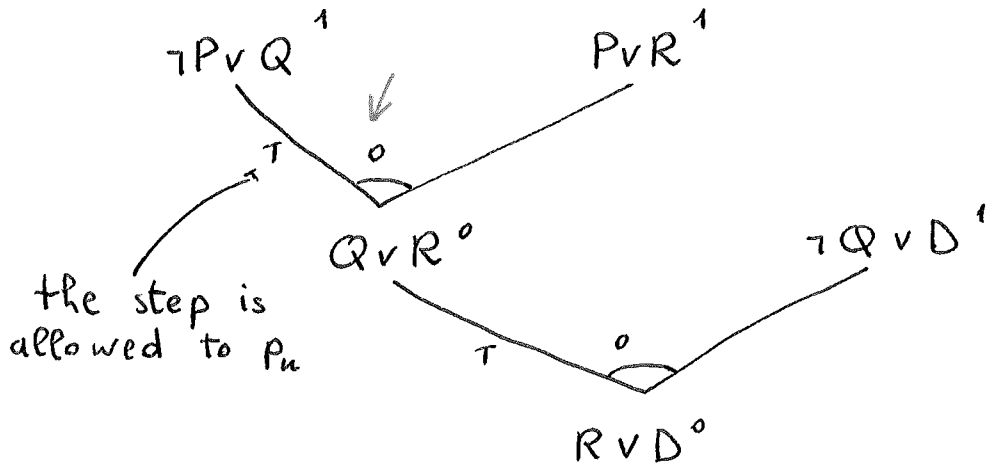$\mathcal{C}'$ :   parallelization by subdivision of $\mathcal{C}$

Is $\mathcal{C}'$ contraction - based ?

## Sufficient conditions :      (two sets)

1)  $\Sigma'$ propagates clauses up to redundancy & does not subdivide contractions

2)  $\Sigma'$ propagates clauses up to redundancy & subdivides generations, not deletions, by contraction

How can we represent distributed search in the search space of a T.P. problem?
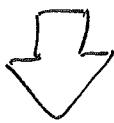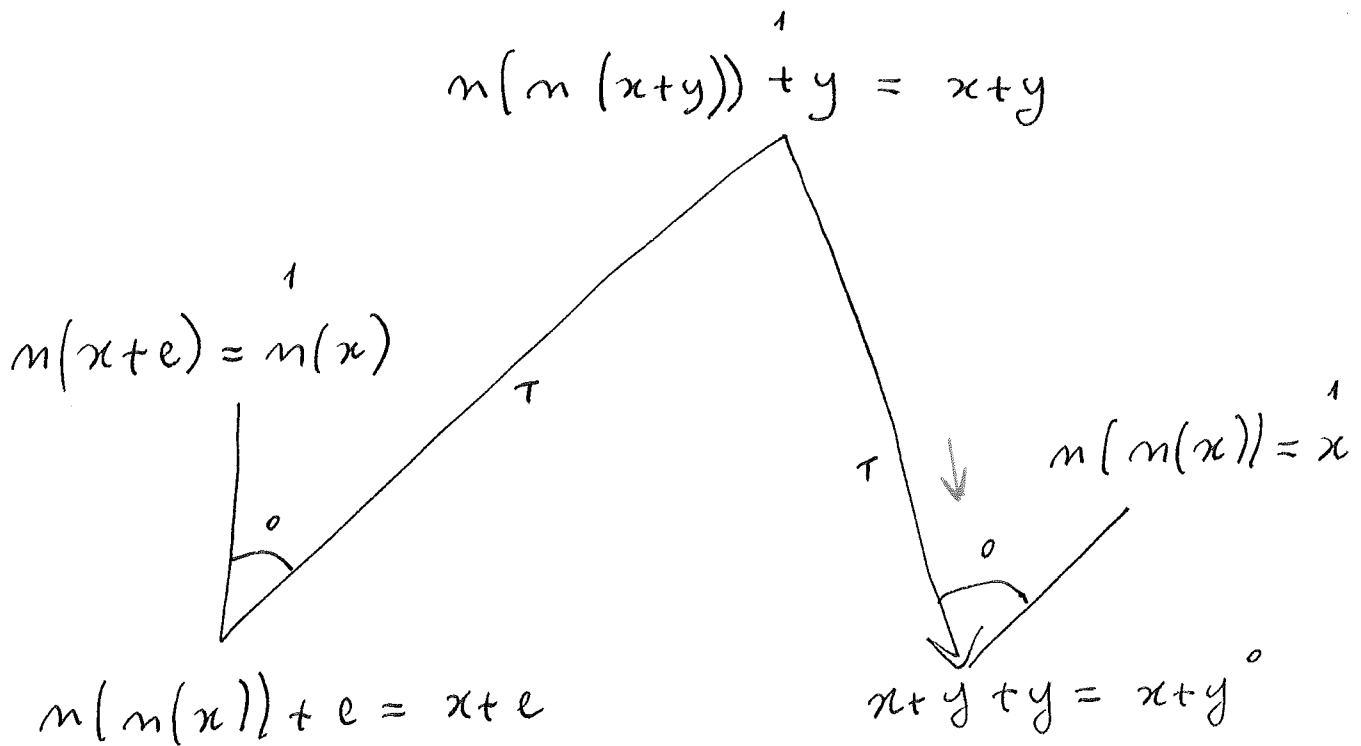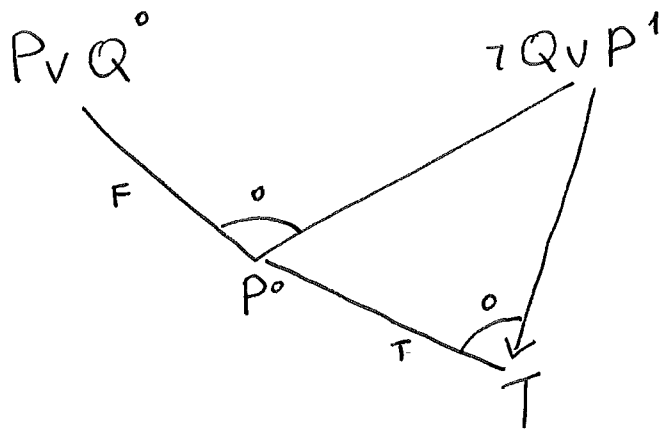
# Example : expansion

$\neg P \vee Q^1$    $P \vee R^1$

$\searrow$

T    o

$\neg T$

the step is
allowed to $P_k$

$Q \vee R^0$    $\neg Q \vee D^1$

T    o

$R \vee D^0$

$\Downarrow$

$\neg P \vee Q^1$    $P \vee R^1$

T    1

$Q \vee R^1$    $\neg Q \vee D^1$

T    o

$R \vee D$

Marking relative to $P_k$ :
have one per process.

# $\underline{Example} : contraction$

$$n\big(n\,(x+y)\big) \overset{1}{+} y \;=\; x+y$$

$$n(x+e) \overset{1}{=} n(x)$$

$T$

$n\big(n(x)\big) + e = x+e$

$\overset{0}{}$

$T$ $\quad\downarrow\quad$ $n\big(n(x)\big) \overset{1}{=} x$

$x+\overset{}{y}+y = x+y\overset{0}{}$

$\Downarrow$

$$n\big(n(x+y)\big) \overset{-1}{+} y = x+y$$

$n(x+e) \overset{1}{=} n(x)$ $\quad T$

$T$ $\qquad n\big(n(x)\big) \overset{1}{=} x$

$n\big(n(x)\big) + e = x+\overset{0}{e}$

$\overset{1}{}$

$x+y+y = x+y$

# Example: communication

$P \lor Q^0$             $\neg Q \lor P^1$

F      0

$P^0$

     0

T     T

at $P_k$

receives $P$
from $P_j$ and
applies it to
subsume $\neg Q \lor P$

$P \lor Q^0$             $\neg Q \lor P^{-1}$

F      1

$P^1$

     1

T     $T^1$

at $P_k$

# Representation of search space

Closure: $S_I^*$

Marked search graph $G(S_I^*) = \langle V, E, \ell, h, \bar{s}, \bar{c} \rangle$

Vertices $V$: clauses $(\ell: V \rightarrow \mathcal{L}/\doteq)$

Hyperarcs $E$: inferences $(h: E \rightarrow I)$

Marking $\bar{s}$ of vertices: $s^n: V \rightarrow \mathbb{Z}$

$s^n(v)$: # of variants of clause
($-1$ if all deleted by $P_n$)

Marking $\bar{c}$ of arcs: $c^n: E \rightarrow \mathbb{N} \times Bool$

$\pi_1(c^n(e)) = $ # of times $P_n$ executed $e$ or received clauses generated by $e$

$\pi_2(c^n(e)) = $ true / false
(allowed / forbidden)

# Evolution of search space

$$S_0^\kappa \vdash \ldots S_i^\kappa \vdash \ldots \qquad \kappa \in [0, n-1]$$

$$e = (v_1 \ldots v_m ; v_{m+1} ; u) \qquad \text{enabled}$$

1) all premises present $\quad (s^\kappa(v) \geqslant 1)$

2) arc allowed $\qquad (\pi_2(c^\kappa(e)) = \text{true})$

$$S_0^\kappa(v) = 0$$

$$S_{i+1}^\kappa(v) = \begin{cases} S_i^\kappa(v) + 1 & \text{if generated or received} \\ \\ S_i^\kappa(v) - 1 & \text{if deleted} \\ (-1 \text{ if last variant}) \end{cases}$$

$$\pi_1(c_0^\kappa(e)) = 0$$

$$\pi_1(c_{i+1}^\kappa(e)) = \pi_1(c_i^\kappa(e)) + 1 \qquad \text{if executed or received}$$

$$\pi_2(c_{i+1}^\kappa(e)) = \begin{cases} \alpha(S_0 \ldots S_{i+1}, m, \kappa, f, \bar{x}) & \text{if} \neq \perp \\ \\ \text{true} & \text{otherwise} \end{cases}$$

How to analyze
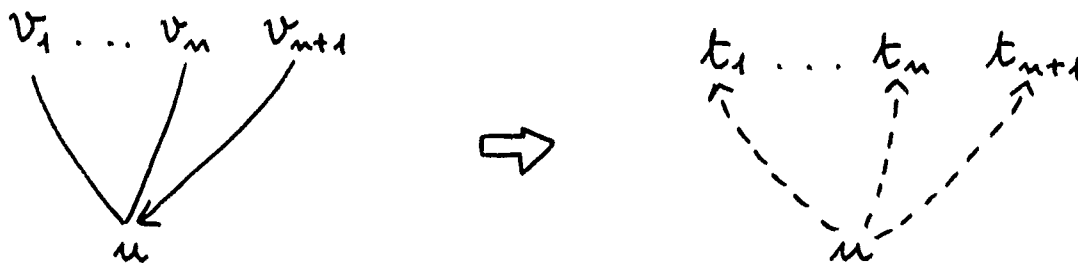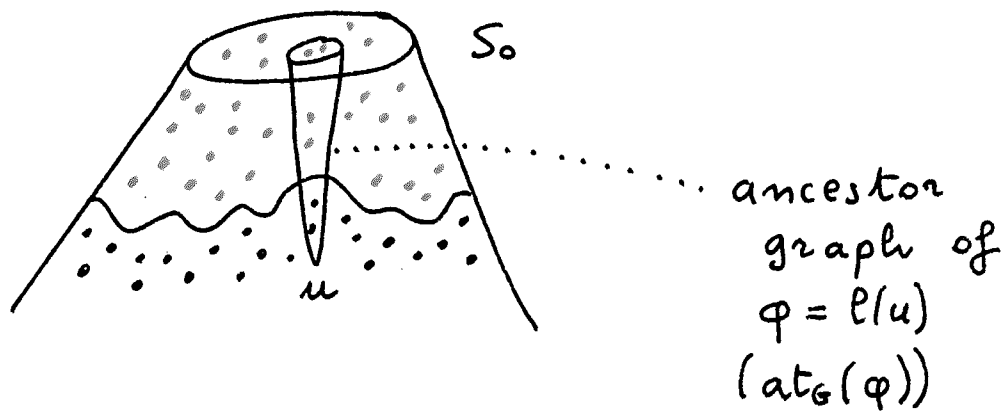
T. P. strategies ?

# Methodological problems

T. P. is only semi-decidable $\Rightarrow$
Search space is infinite
algorithm analysis does not apply.

Complexity proportional
neither to input (e.g., input length)
nor to output (e.g., proof length).

Need to analyze the <u>search process</u>:
for all factors ( communication,
overlap, parallel searches, subdivision)
find suitable representation and
measure benefit / cost.

# Measuring search complexity

G:



$S_0$

ancestor graph of $\varphi = \ell(u)$ $(at_G(\varphi))$



$$t = (u \, ; \, e \, ; \, (t_1 \ldots t_{m+1}))$$

where $t_i$ ancestor-graph of $v_i$

we $t$ is relevant to $u$ in $t$ for $P_k$ if

- $w \in \{v_1 \ldots v_{m+1}\}$ and $\pi_1(c^k(e)) = 0$ or

- $w$ is relevant to $v_i$ in $t_i$ for some $i$

# A notion of distance in search spaces

Past distance:

$$pdist_{G^\kappa}(t) = |\{w|\ w \in t,\ \overset{*}{s}(w) \neq 0\}|$$

Future distance:

$$fdist_{G^\kappa}(t) = \begin{cases} \infty & \text{if } \overset{*}{s}(\varphi) < 0 \text{ or} \\ & \exists w \in Rev_{G^\kappa}(t)\ \overset{*}{s}(w) < 0 \\ |\{w|\ w \in t,\ s(w) = 0\}| \end{cases}$$

Global distance:

$$gdist_{G^\kappa}(t) = pdist_{G^\kappa}(t) + fdist_{G^\kappa}(t)$$

$$fdist_{G^\kappa}(\varphi) = \min\ fdist_{G^\kappa}(t) \qquad t \in at_G(\varphi)$$

Dynamic distance:

$fdist_{G^\kappa}(t)$ measures the part of $t$ that $P_\kappa$ needs to traverse to reach $\varphi$ via $t$

if $\infty$, then unreachable (redundant)

# Bounded search spaces

At stage $i$ ($i \geqslant 0$) of a derivation
define the _bounded search space_
reachable (by process $P_\kappa$) within
distance $j$ ($j > 0$) from the start:

$$
\text{space}(G^\kappa, j) = \sum_{\substack{v \in V \\ v \neq T}} \text{mul}_{G^\kappa}(v, j) \cdot \ell(v)
$$

where

$$
\text{mul}_{G^\kappa}(v, j) = \left| \left\{ t : \begin{array}{l} t \in at_G(v), \\ t \text{ allowed for } P_\kappa \\ 0 < gdist_{G^\kappa}(t) \leqslant j \end{array} \right\} \right|
$$

Ancestor - graph _forbidden_ for $P_\kappa$ if
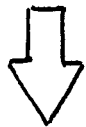$\exists e \ \ \pi_1(c^\kappa(e)) = 0$ and $\pi_2(c^\kappa(e)) = $ false,
_allowed_ otherwise.

# Analysis of the search process in distributed search

# Subdivision

Ancestor graph forbidden for $P_n$ if
$\exists e \quad \pi_1(c^n(e)) = 0 \quad$ and $\quad \pi_2(c^n(e)) = false$,
allowed otherwise.

$$\underset{P^0}{\bigvee\,}_0 \Big/ T$$

$$\bigvee\,_0 \Big/ T$$

allowed

$$\Downarrow \qquad \text{generates or receives } P$$
$$(\alpha \text{ becomes defined}_{\text{on } P})$$

$$\underset{P^1}{\bigvee\,}_1 \Big/ T$$

$$\bigvee\,_0 \Big/ T$$

allowed

$$\underset{P^1}{\bigvee\,}_1 \Big/ T$$

$$\bigvee\,_0 \Big/ F$$

forbidden

# Communication

$$\bigvee \!\!{}^{0}\!/F$$

$$P^{0}$$

$$\bigvee \!\!{}^{0}\!/T$$

forbidden

⇩    receives P  ($\alpha$ becomes defined)

$\bigvee \!\!{}^{1}\!/F$      $\bigvee \!\!{}^{1}\!/F$

$P^{1}$        $P^{1}$

$\bigvee \!\!{}^{0}\!/T$      $\bigvee \!\!{}^{0}\!/F$

allowed       forbidden

$P_n$ and $P_h$ overlap on $t$ if $t$ is allowed for both

# Contraction and communication

Sequential contraction - based:

if a deleted $\varphi$ is re-generated, it is deleted again (monotonic) before being used (eager).

$$\Downarrow$$

$$\oint dist_i(t) = \infty \quad \text{then} \quad \forall j > i \quad \oint dist_j(t) = \infty$$

In parallel:

assume local eager contraction:

if a deleted or unreachable $\varphi$ is re-generated or received, it is deleted again before being used.

$$\Downarrow$$

$$\oint dist_{G_i^k}(\varphi) = \infty \quad \text{then} \quad \forall j \quad \oint dist_{G_j^k}(\varphi) = \infty$$

# Contraction and communication

$\boxed{P_\kappa}$

$\overset{-1}{Q}$

$R^0$ (not redundant)

$\oint dist(t) = \infty$

receives $R$ $\longleftarrow$

$\overset{-1}{Q}$

$R$

$\oint dist(t) \neq \infty$

$\boxed{P_\ell}$

$\overset{1}{Q}$

$R^0$

generates $R$ from $Q$ and sends it

$Q$ still deleted but no longer relevant

$P_\ell$ : late contraction

$P_\kappa$ : contraction undone

# Evolution of bounded search spaces

1) if $S_i^{\kappa} \vdash S_{i+1}^{\kappa}$ generates $\psi$

   $\forall j \quad space(G_{i+1}^{\kappa}, j) \preceq_{mul} space(G_i^{\kappa}, j)$

   because of subdivision

2) if $S_i^{\kappa} \vdash S_{i+1}^{\kappa}$ replaces $\psi$ by $\psi'$

   $\forall j \quad space(G_{i+1}^{\kappa}, j) \preceq_{mul} space(G_i^{\kappa}, j)$

   because of contraction and
   
   subdivision

3) if $S_i^{\kappa} \vdash S_{i+1}^{\kappa}$ receives $\psi$

   $\forall j \quad \exists l \preceq i \quad space(G_{i+1}^{\kappa}, j) \preceq_{mul} space(G_l^{\kappa}, j)$

   because of subdivision,
   subdivision undone and contraction
   undone

$\Downarrow$

*Now - monotonic* bounded search spaces

## Parallel bounded search spaces

$$gmul_G(v, j) = \sum_{k} mul_{G^k}(v, j)$$
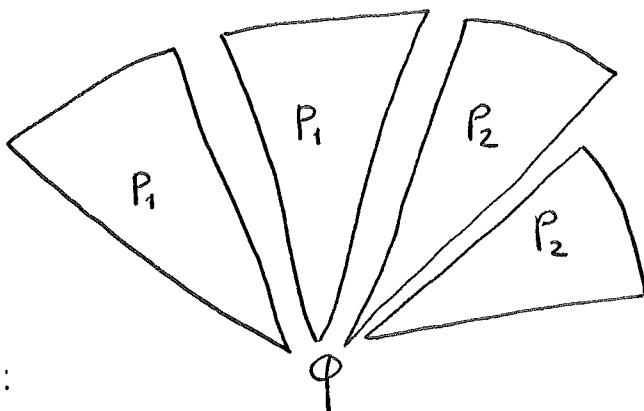
$$pmul_G(v, j) = \lfloor gmul_G(v, j) / w \rfloor$$

where $w$ is the $\#$ of processes

$$pspace(G, j) = \sum_{\substack{v \in V \\ v \notin T}} pmul_G(v, j) \cdot \ell(v)$$
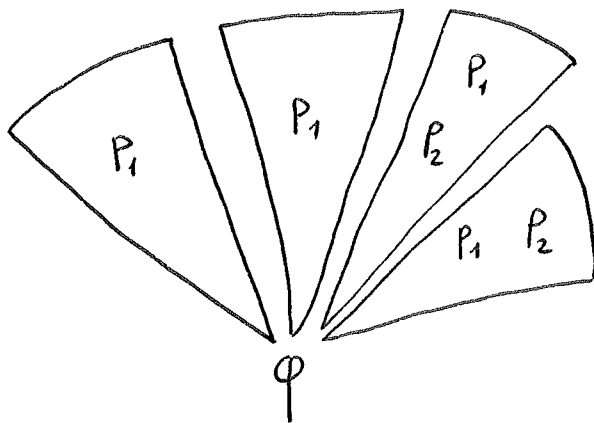
Capture overlap of the processes

# Example

$P_1$  $P_2$



No overlap:

$\text{mul}_{G^1}(\varphi, j) = 2$ $\qquad$ $\text{mul}_{G^2}(\varphi, j) = 2$

$\text{gmul}_G(\varphi, j) = 4$ $\qquad$ $\text{pmul}_G(\varphi, j) = 2$



Overlap:

$\text{mul}_{G^1}(\varphi, j) = 4$ $\qquad$ $\text{mul}_{G^2}(\varphi, j) = 2$

$\text{gmul}_G(\varphi, j) = 6$ $\qquad$ $\text{pmul}_G(\varphi, j) = 3$

"average"

# Minimize overlap

1) Overlap due to inaccurate subdivision

2) Overlap due to communication

Two properties of $\alpha$:

No arc-duplication :   avoid (1)

No clause-duplication:   minimize (2)

Lemma:   local eager contraction +

   no clause-duplication $\Rightarrow$

$P_k$: allowed to generate $\varphi$

$\exists r \quad \forall k \neq h \quad \forall i \geqslant r \quad \forall j$

$$mul_{G_i^k}(\varphi, j) \leqslant 1$$

How to compare a distributed - search contraction - based strategy with its sequential base ?

# Analysis of $\mathcal{C}$ vs. $\mathcal{C}'$

$\mathcal{C} = \langle I, \leq \rangle$ sequential contraction-based

$\mathcal{C}' = \langle I, \leq' \rangle$ contraction-based parallelization
by subdivision of $\mathcal{C}$

Same $I$ $\Rightarrow$ same initial search space

Lemmas:

1) $\varphi \in S_i \Rightarrow \exists_{P_k} \exists_j \quad \varphi \in S_j^k \cup R(S_j^k)$

2) $\varphi \in R(S_i) \Rightarrow \forall_{P_k} \exists_j \quad \varphi \in R(S_j^k)$

3) $fdist_{G_i}(\varphi) = \infty \Rightarrow \forall_{P_k} \exists_j \forall \ell \geqslant j$

    either $fdist_{G_\ell^k}(\varphi) = \infty$

    or $\forall t \in at_G(\varphi) \quad t$ forbidden

$$\Downarrow$$

Theorem:

$$fdist_{G_i}(\varphi) = \infty \Rightarrow \exists_r \forall i \geqslant r \forall_j$$
$$pmul_{G_i}(\varphi, j) = 0$$

# A limit lemma

Local eager contraction +
immediate propagation
(hence global eager contraction)

$$\psi' \in S_\infty - R(S_\infty)$$

if $s_i^u(\psi) = -1$ $\qquad \psi \in Rev_{G_i^u}(t)$ $\qquad$ for $P_u$ then:

1) $\forall P_u \ \forall j \quad s_j^u(\psi) = -1 \quad \Rightarrow \quad \psi \in Rev_{G_j^u}(t)$

(what is relevant for a process is relevant
for all : no late contraction)

2) $\forall j \geqslant i \quad \psi \in Rev_{G_j^u}(t)$

(what is relevant remains relevant:
no contraction undone)

$\Downarrow$

$f dist_i(t) = \infty$ $\qquad$ then $\qquad \forall j > i \quad f dist_j(t) = \infty$

# A limit theorem

Assume:

immediate propagation of clauses up to redundancy

no clause - duplication

Lemma:

$$\text{dist}_{G_i}(\varphi) \neq \infty \quad \forall i \implies \exists z \; \forall i \geqslant z \; \forall j$$
$$\text{pmul}_{G_i}(\varphi,j) \leqslant \text{mul}_{G_i}(\varphi,j)$$

Theorem:

$$\forall j \; \exists m \; \forall i \geqslant m \quad \text{pspace}(G_i,j) \leqslant_{\text{mul}} \text{space}(G_i,j)$$

Significance:

1) "limit theorem" that strategies may approximate (e.g., by reducing overlap)

2) "negative" result which contributes to explain intrinsic difficulty of parallel theorem proving

# Discussion

Strategy analysis: study of search
in infinite search spaces

Model: marked search graph
Measure: bounded search spaces
Already applied to analysis of contraction
Now: distributed search

## Analytic comparison:

"Limit theorem" explains nature of problem
(overlap + communication/contraction)

When adopting asynchronous distributed search
one expects that contraction may be delayed,
but synchronizing on every inference is hopeless,
and one may conjecture subdivision compensates
for late contraction: not so in general
<div align="right">(worst-case scenario).</div>

Relevant to problems where eager contraction
is important: not a small class based on
experience.

# Directions for future work

On analysis of parallel search:

- Reordering of search
  relevant to both
  distributed search
  multi - search

On analysis of theorem proving:

- Comparison of search plans
- Subgoal - reduction strategies
- Reasoning modulo a theory

# Distributed search for CBS

## Clause – Diffusion (1992)

| | | |
|---|---|---|
| Aquarius | (Otter 2.2) | 1992 |
| Peers | (OPS 1/93) | 1993-94 |

## Modified Clause – Diffusion (1994-96)

| | | |
|---|---|---|
| Peers-mcd | (EQP 0.9) | 1996-98 |
| " | (EQP 0.9d) | 1999 |
| " | + hybrid mode | 2000 |

Hybrid mode: distributed search + multi-search

Levi Commutator Problem in group theory: super-linear speed-up

Robbins algebras are Boolean: super-linear speed-up

Moufang identities in alternative rings without cancellation laws built-in (EQP cannot do)