

PAI

PETR-MAREK.COM

- DOMÉNOVÉ NEZÁVISLÉ PLÁNOVÁNÍ
- VLASTNOSTI, HEURISTIKY A ALGORITMY

- PLANNING PROBLEM

- STATES
- INITIAL STATE
- GOAL STATE
- ACTIONS THAT TRANSFORMS ONE STATE TO ANOTHER

- PLAN

- SEQUENCE OF ACTIONS THAT TRANSFORMS INITIAL STATE TO GOAL STATE

- POSSIBLE TASKS

- FIND OUT WHETHER THERE IS A SOLUTION
- FIND ANY SOLUTION
- FIND OPTIMAL SOLUTION
- FIND BEST SOLUTION IN GIVEN TIME
- FIND SOLUTION SATISFYING ~~BY~~ SOME PROPERTY

- STATE MODEL FOR CLASSICAL AI PLANNING

- FINITE SPACE STATE  $S$
- INITIAL STATE  $s_0 \in S$
- SET OF GOAL STATES  $S_G \subseteq S$
- APPLICABLE ACTIONS  $A(s) \subseteq A$  FOR  $s \in S$
- TRANSITION FUNCTION  $s' = f(a, s)$  FOR  $a \in A(s)$
- COST FUNCTION  $c: A^* \rightarrow [0, \infty)$

## - TRANSITION SYSTEM

$$- \langle S, I, \{\alpha_1, \dots, \alpha_n\}, G \rangle$$

-  $S$  - FINITE SET OF STATES

-  $I \subseteq S$  IS A FINITE SET OF INITIAL STATES

- ACTION  $\alpha_i \subseteq S \times S$  IS BINARY RELATION ON  $S$

-  $G \subseteq S$  IS FINITE SET OF GOAL STATES

## - APPLICABLE ACTION

- ACTION  $\alpha$  IS APPLICABLE IN STATE  $s$  IF  $s \alpha s'$  EXIST FOR AT LEAST ONE STATE  $s'$

## - DETERMINISTIC TRANSITION SYSTEM

- ONLY ONE INITIAL STATE

- ALL ACTIONS ARE DETERMINISTIC

- HENCE ALL FUTURE STATES ARE COMPLETELY PREDICTABLE

$$- \langle S, I, O, G \rangle$$

-  $S$  IS FINITE SET OF STATES

-  $I \subseteq S$  IS INITIAL STATE

- ACTIONS  $\alpha \in O$  (WITH  $\alpha \subseteq S \times S$ ) ARE PARTIAL FUNCTIONS

-  $G \subseteq S$  IS A FINITE SET OF GOAL STATES

## - SUCCESSOR STATE

$$- s' = \text{APP}_\alpha(s)$$

## - PLAN

- FOR  $\langle S, I, A, G \rangle$  IS SEQUENCE  $\pi$

$$\pi = a_1 \dots a_m$$

- SUCH THAT

$a_1, \dots, a_m \in A$  AND  $s_0, \dots, s_m$  IS SEQUENCE OF STATES

-1.  $s_0 = I$

-2.  $s_i = \text{APP}_{a_i}(s_{i-1})$  FOR EVERY  $i \in \{1, \dots, m\}$

-3.  $s_m \in G$

- THIS CAN BE ALSO EXPRESSED AS

$$\text{APP}_{a_m}(\text{APP}_{a_{m-1}}(\dots \text{APP}_{a_1}(I))) \in G$$

## - DETERMINISTIC PLANNING TASK

- 4-TUPLE  $\langle V, I, A, G \rangle$

-  $V$  IS FINITE SET OF STATE VARIABLES

-  $I$  IS INITIAL STATE OVER  $V$

-  $A$  IS FINITE SET OF ACTIONS OVER  $V$

-  $G$  IS CONSTRAINT OVER  $V$  DESCRIBING THE GOAL STATES

## - SAS

- TUPLE  $\langle V, A, I, G \rangle$

-  $V$  IS FINITE SET OF STATE VARIABLES WITH FINITE DOMAINS  $\text{dom}(v_i)$

-  $I$  IS INITIAL STATE OVER  $V$

-  $G$  IS PARTIAL ASSIGNMENT TO  $V$

-  $A$  IS FINITE SET OF ACTIONS  $\alpha$  SPECIFIED BY  $\text{PRE}(\alpha)$  AND  $\text{EFF}(\alpha)$  BOTH BEING PARTIAL ASSIGNMENTS TO  $V$

## - STRIPS

- TUPLE  $\langle P, A, I, G \rangle$

-  $P$  IS FINITE SET OF ATOMS (BOOLEAN VARS)

-  $I \in P$  IS INITIAL SITUATION

-  $G \in P$  IS GOAL SITUATION

-  $A$  IS A FINITE SET OF ACTIONS  $a$  SPECIFIED BY  $PRE(a)$ ,  
 $ADD(a)$  AND  $DEL(a)$ , ALL SUBSETS OF  $P$

## - PROGRESSION PLANNERS

- FIND SOLUTION BY FORWARD SEARCH

- SEARCH SPACES

- SEARCH NODES CORRESPOND TO STATES

- WHEN WE ENTER THE SAME STATE ALONG DIFFERENT PATH  
WE DON'T HAVE TO CONSIDER THE STATE AGAIN

- SEARCH NODES CORRESPOND TO OPERATOR SEQUENCES

- BUT DIFFERENT SEQUENCES CAN LEAD TO IDENTICAL STATES

## - BACKWARD VS. FORWARD SEARCH

- IT IS NOT SYMMETRIC

- FORWARD

- ONE INITIAL TO SET OF GOALS

- WHEN WE APPLY OPERATOR  $o$  IN STATE, THERE IS SINGLE UNIQUE  
SUCCESSOR  $s'$

- BACKWARD

- SEVERAL GOAL STATES TO ONE INITIAL STATE

- APPLICATION OF OPERATOR  $o$  IN  $s'$  CAN LEAD TO SEVERAL  $s$

## - CLASSIFICATION OF SEARCH ALGORITHMS

- UNIFORMED

- INFORMED (HEURISTIC)

- SYSTEMATIC ALGORITHMS

- CONSIDER LARGE NUMBER OF SEARCH NODES SIMULTANEOUSLY

- LOCAL SEARCH ALGORITHMS

- WORK WITH FEW CANDIDATE SOLUTIONS AT A TIME

- ~~OF~~ USEFULNESS

- SATISFYING PLANNING

- HEURISTIC SEARCH OUTPERFORMS

- AND SOMETHING BETWEEN SYSTEMATIC AND LOCAL SEARCH

- OPTIMAL

- UNIFORMED IS BETTER

- SYSTEMATIC IS REQUIRED

- UNIFORMED, SYSTEMATIC

- DFS, BFS, ITERATIVE-DFS

- UNIFORMED, LOCAL

- RANDOM WALK

- ~~SYSTEMATIC~~ HEURISTIC, SYSTEMATIC

- GREEDY BEST-FIRST

-  $A^*$

- WEIGHTED  $A^*$

- IDA $^*$

## - HEURISTIC LOCAL

- HILL CLIMBING

- BEAM SEARCH

- GENETIC ALGORITHMS

- SIMULATED ANNEALING

## - HEURISTIC SEARCH

### - HEURISTIC FUNCTION

$$- h: \Sigma \rightarrow \mathbb{N}_0 \cup \{\infty\}$$

- ESTIMATES DISTANCE FROM  $\sigma$  (STATE) TO NEAREST GOAL NODE

- EFFICIENCY OF ALGORITHM CLOSELY RELATES TO HOW ACCURATELY  $h$  REFLECTS THE ACTUAL GOAL DISTANCE

### - OPTIMAL (PERFECT) HEURISTIC

$$- h^*$$

- MAPS EACH SEARCH NODE TO LENGTH OF SHORTEST PATH TO ANY GOAL STATE

### - SAFE

$$- h^*(\sigma) = \infty \text{ FOR ALL } \sigma \in \Sigma \text{ WITH } h(\sigma) = \infty$$

- " RETURNS INFINITY FOR NODES WITHOUT CONNECTION TO GOAL STATES "

### - GOAL-AWARE

$$- h(\sigma) = 0 \text{ FOR ALL GOAL NODES } \sigma \in \Sigma$$

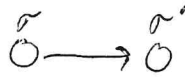
### - ADMISSIBLE

$$- h(\sigma) \leq h^*(\sigma)$$

- " ALWAYS UNDERESTIMATE "

- CONSISTENT

$$- h(\sigma) \leq h(\sigma') + 1$$



- "HEURISTIC IS DECREASED BY COST OF PATH IF WE TRAVEL CLOSER TO GOAL"

- WEIGHTED  $A^*$

$$- f(\sigma) = g(\sigma) + W \cdot h(\sigma)$$

- IF  $W = 1 \dots A^*$

- IF  $W = 0 \dots$  BREATH FIRST (DIJKSTRA)

- IF  $W \rightarrow \infty \dots$  GREEDY BEST FIRST

- GENERAL IDEA OF ADMISSIBLE HEURISTIC

- OBTAINED AS OPTIMAL COST FUNCTION OF RELAXED PROBLEMS

- FOR EXAMPLE

- EUCLIDIAN DISTANCE

- MANHATTAN DISTANCE

- SPANNING TREE IN TSP

- DOMINATION OF HEURISTICS

- IF  $h_2(\sigma) \geq h_1(\sigma)$  FOR ALL NODES  $\sigma$

- THEN  $h_2$  DOMINATES  $h_1$

-  $h_2$  IS BETTER FOR SEARCH

- COMBINING ADMISSIBLE HEURISTICS

- ADMISSIBLE HEURISTICS  $h_1, \dots, h_n$

$$- h(\sigma) = \max_{i=1}^n \{h_i(\sigma)\}$$

- IT IS ALSO ADMISSIBLE HEURISTIC AND DOMINATES ALL INDIVIDUAL  $h_i$

## - HEURISTIC USED IN STRIPS

- NUMBER OF STATE ~~HEURISTIC~~ VARIABLES THAT DIFFER IN CURRENT STATE AND STRIPS GOAL

- "MORE TRUE GOAL LITERALS  $\rightarrow$  CLOSER TO GOAL"

- BUT NOT THAT GREAT

- NOT INFORMATIVE

- SMALL RANGE OF VALUES

- MANY SUCCESSORS HAVE THE SAME VALUE

- SENSITIVE TO REFORMULATION OF TASKS

- IGNORES PROBLEM STRUCTURE

- VALUE DOESN'T DEPEND ON SET OF ACTIONS

- GENERAL PROCEDURES TO COME UP WITH HEURISTICS

- "SOLVE EASIER VERSION OF PROBLEM"

- RELAXATION

- LESS CONSTRAINED PROBLEM

- ABSTRACTION

- SMALLER VERSION OF PROBLEMS



## - RELAXATION

- "WE WOULD LIKE TO IGNORE BAD SIDE EFFECTS OF ACTIONS"

### - EXAMPLE

- 8-PUZZLE

- IF WE MOVE TILE FROM X TO Y, GOOD EFFECT IS THAT X IS NOW FREE

- BAD EFFECT IS THAT Y IS NOT FREE, PREVENTING US TO MOVE TILES THROUGH IT

### - STRIPS

- EFFECTS MAKING ATOMS TRUE

- GOOD

- EFFECTS MAKING ATOMS FALSE

- BAD

- IGNORE ALL DELETE EFFECTS

### - RELAXATION:

$$\alpha = \langle \text{PRE}(\alpha), \text{ADD}(\alpha), \text{DEL}(\alpha) \rangle \rightsquigarrow \alpha' = \langle \text{PRE}(\alpha), \text{ADD}(\alpha), \emptyset \rangle$$

- CAN BE SOLVED BY GREEDY ALGORITHM

- IT IS SOUND

- IF RETURNS A PLAN, IT IS CORRECT SOLUTION

~~IF~~ - IF RETURNS UNSOLVABLE, THEN PROBLEM IS UNSOLVABLE

- RUNS IN POLYNOMIAL TIME  $O(|P|)$

- POSSIBLE USABILITY OF RELAXATION

- USE OPTIMAL PLANNER FOR RELAXED PROBLEM

- NP-HARD

-  $R^+$  HEURISTIC

- COMPUTE ESTIMATION OF ITS DIFFICULT IN DIFFERENT WAY

-  $R_{MAX}$

-  $R_{ADD}$

- COMPLETE SOLUTION NOT OPTIMAL, BUT REASONABLE

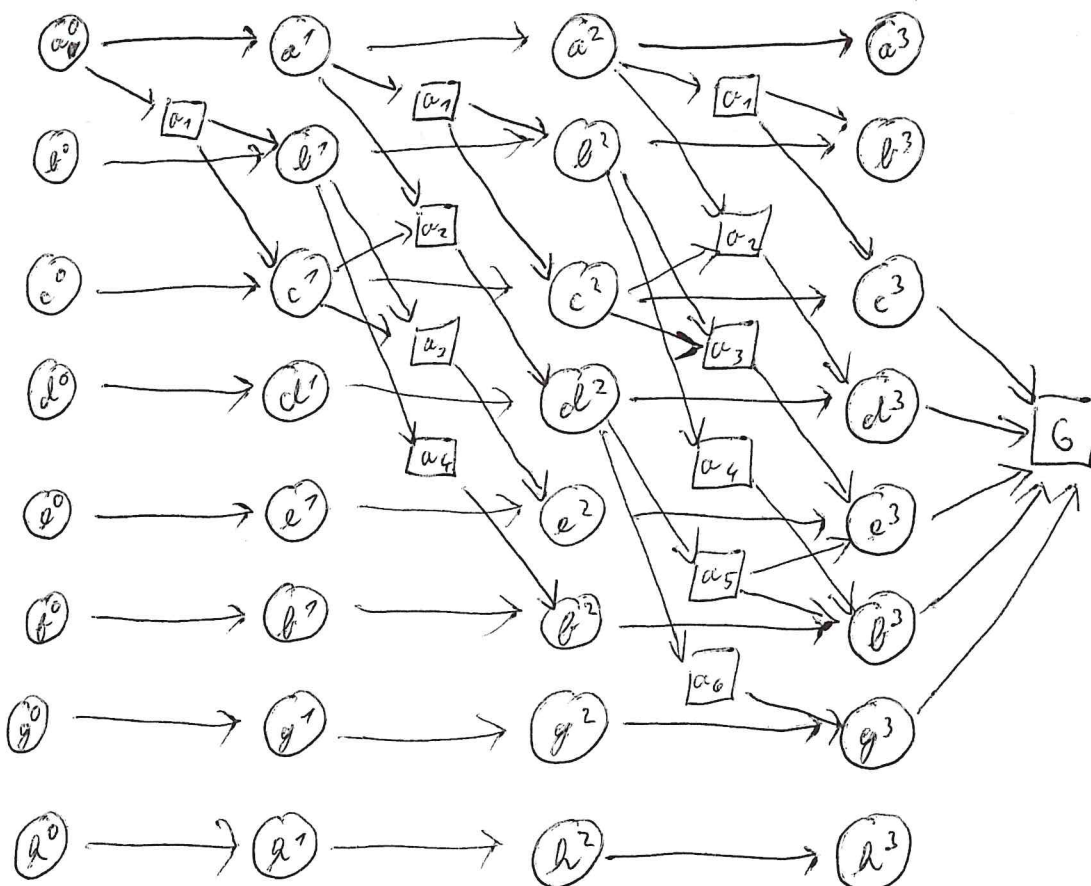
-  $R_{FF}$

- GRAPHICAL INTERPRETATION

$$I = \{ a=1, b=0, c=0, d=0, e=0, f=0, g=0, h=0 \}$$

$$a_1 = \langle \{a\}, \{b, c\}, \emptyset \rangle \quad a_2 = \langle \{a, c\}, \{d\}, \emptyset \rangle \quad a_3 = \langle \{b, c\}, \{e\}, \emptyset \rangle \quad a_4 = \langle \{e\}, \{f\}, \emptyset \rangle$$

$$a_5 = \langle \{d\}, \{g\}, \emptyset \rangle$$



- SIMPLEST RELAXED PLANNING HEURISTICS ARE FORWARD COST HEURISTICS

-  $h_{MAX}$ ,  $h_{ADD}$

- NODE APPROXIMATORS ARE COST VALUES

- COST OF NODE IS ESTIMATE HOW EXPENSIVE IN TERMS OF REQUIRED OPERATORS IS TO MAKE THIS NODE TRUE

- PROPAGATE COST BOTTOM UP

- USE COMBINATION RULE FOR

- ACTION NODE

- PROPOSITIONAL NODE

- AT ACTION NODE ADD 1 AFTER COMBINATION RULE

- HEURISTIC VALUE IS COST OF AUXILIARY GOAL NODE

-  $h_{MAX}$

- ACTION NODES

$$\text{COST}(n) = \text{MAX}(\{\text{COST}(n_1), \dots, \text{COST}(n_k)\}) + 1$$

$$\text{MAX}(\emptyset) = 0$$

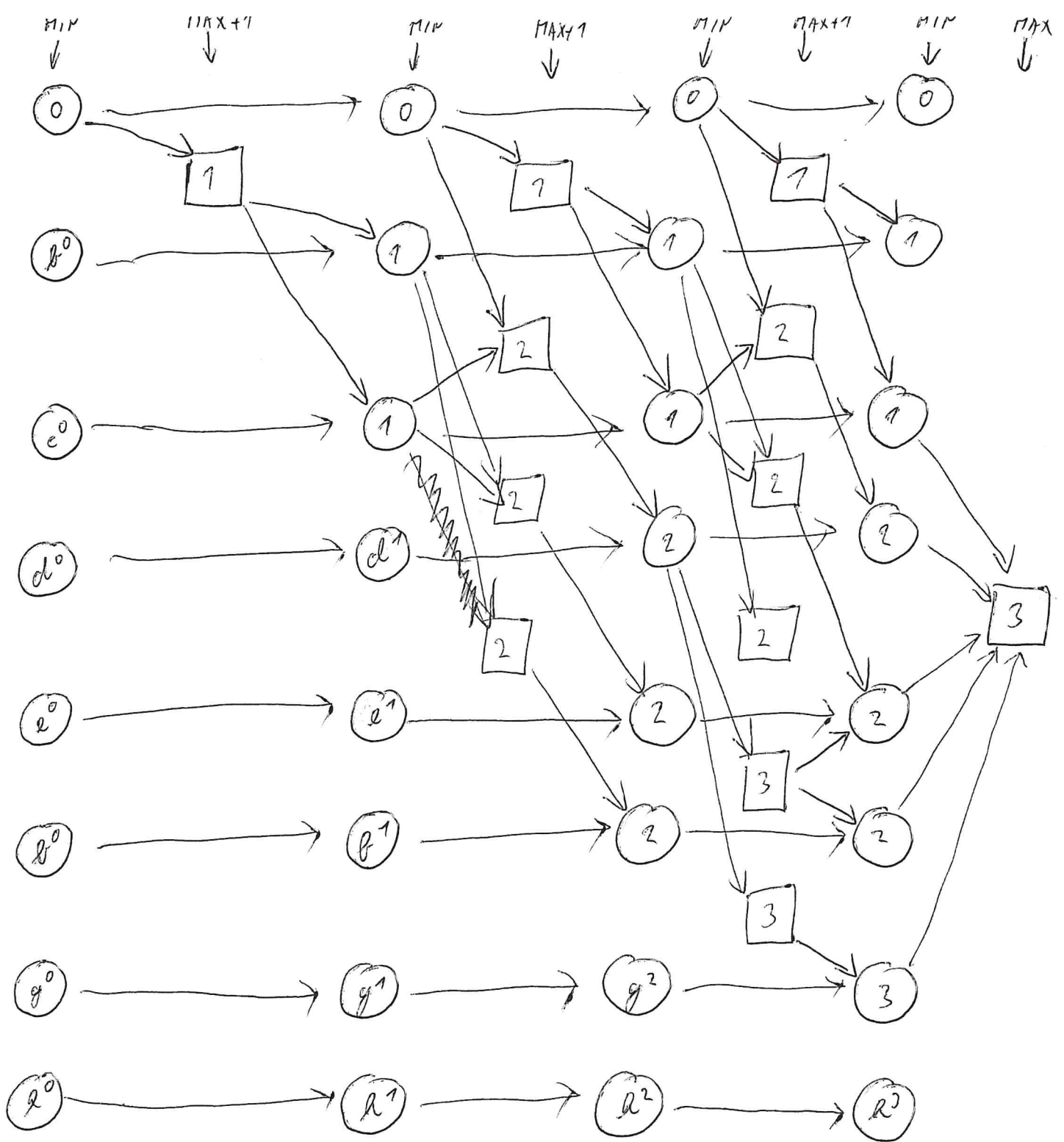
- PROPOSITIONAL

- IF WE HAVE TO ACHIEVE SEVERAL PRECONDITION, ESTIMATE COST BY THE MOST EXPENSIVE ONE

- PROPOSITIONAL NODES

$$\text{COST}(n) = \text{MIN}(\{\text{COST}(n_1), \dots, \text{COST}(n_k)\})$$

- IF WE HAVE CHOICE OF WAY HOW TO ACHIEVE PROPOSITION, SELECT THE CHEAPEST ONE



-  $R_{ADD}$

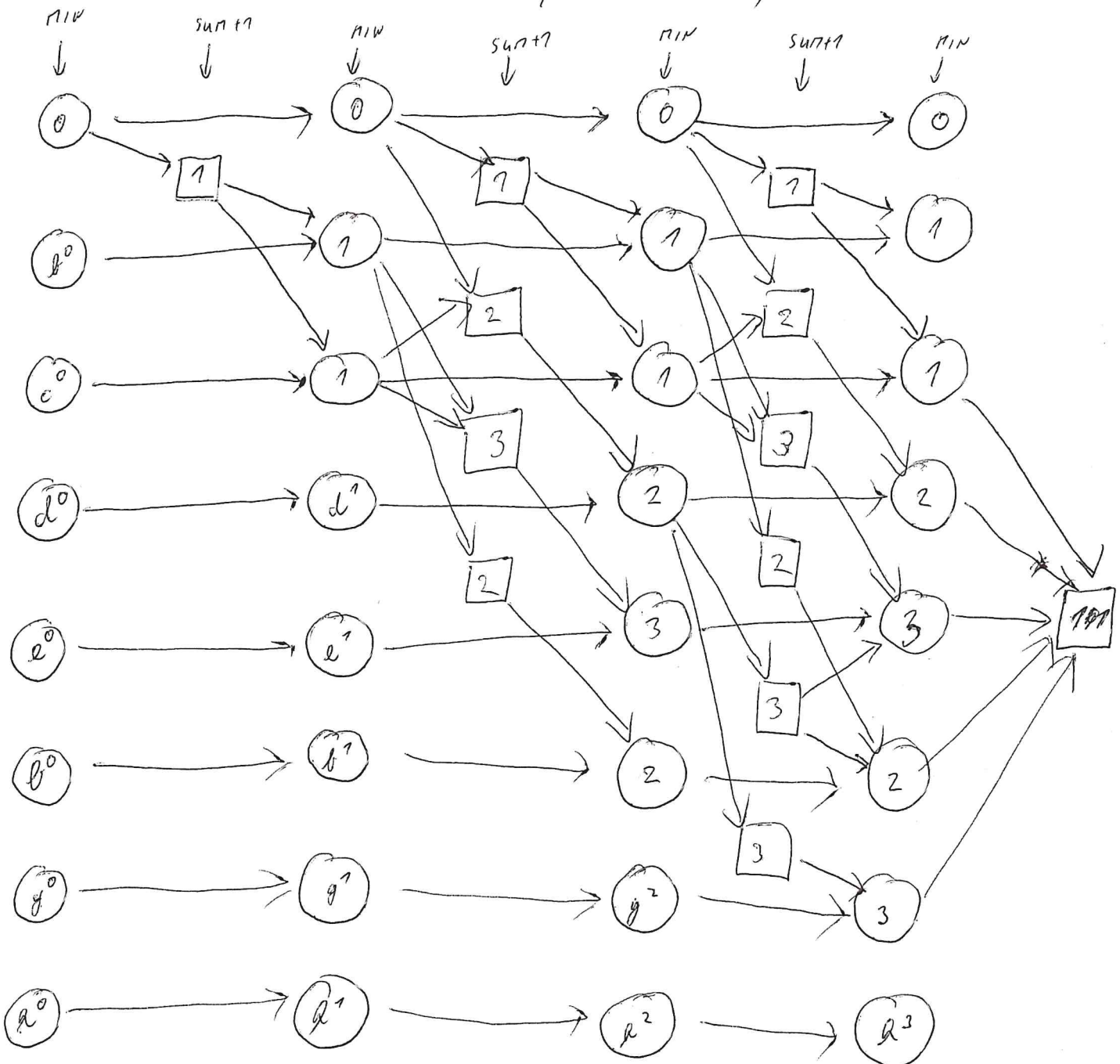
- ACTION MODES

$$COST(n) = (COST(n_1) + \dots + COST(n_k)) + 1$$

- IF WE HAVE TO ACHIEVE SEVERAL PRECONDITIONS, ESTIMATE IT BY COST OF ACHIEVING EACH IN ISOLATION

- PROPOSITIONAL MODES

$$COST(n) = \text{MIN} \{ \{ COST(n_1), \dots, COST(n_k) \} \}$$



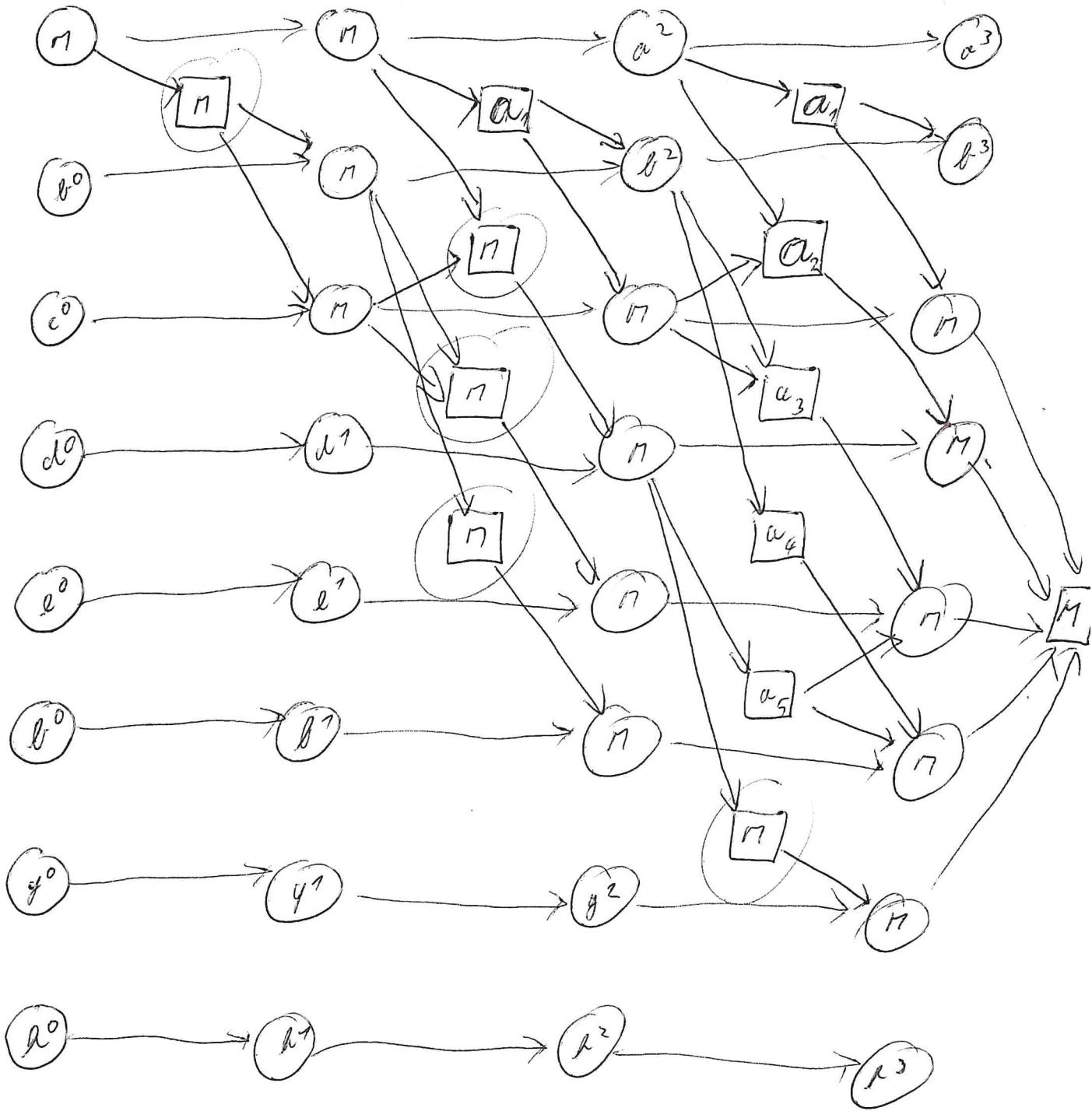
- SAFE, GOAL-AWARE

- VERY INFORMATIVE

- NOT ADMISSIBLE AND, THUS NOT FOR OPTIMAL PLANNING

-  $h_{FF}$

- ANNOTATIONS ARE BOOLEAN VARIABLES
- COMPUTED TOP-DOWN
- GOAL NODE IS INITIALLY MARKED
- JUSTIFIED ACTION MODE
  - ALL TRUE IMMEDIATE PREDECESSORS ARE MARKED
- JUSTIFIED PROPOSITION MODE
  - AT LEAST ONE IMMEDIATE PREDECESSOR IS MARKED
- FOUR RULES TO ANNOTATE
  - MARK ALL IMMEDIATE PREDECESSORS OF MARKED UNJUSTIFIED ACTION MODE
  - MARK IMMEDIATE PREDECESSOR OF MARKED UNJUSTIFIED PROP NODE WITH ONLY ONE IMMEDIATE PREDECESSOR
  - MARK IMMEDIATE PREDECESSOR OF MARKED UNJUSTIFIED PROP NODE CONNECTED VIA 100% ARC
  - MARK ANY IMMEDIATE PREDECESSOR OF MARKED UNJUSTIFIED PROP NODE
- EARLIER RULES ARE MORE PREFERRED
- VALUE OF HEURISTIC IS NUMBER OF MARKED ACTION MODES



$$R_{FF} = 5$$

- SAFE, GOAL AWARE
- NOT ADMISSIBLE, NOT CONSISTENT
- ALWAYS MORE ACCURATE THAN  $R_{ADD}$
- CAN BE COMPUTED IN LINEAR TIME

## - RELATION BETWEEN RELAXATION HEURISTICS

$$h_{\text{MAX}}(s) \leq h_{\text{ADD}}^+(s) \leq h^*(s) \quad \begin{array}{l} \text{OPTIMAL ON WHOLE} \\ \text{PROBLEM} \end{array}$$

$$h_{\text{MAX}}(s) \leq h_{\text{ADD}}^+(s) \leq h_{\text{FF}}(s) \leq h_{\text{ADD}}(s) \quad \begin{array}{l} \text{OPTIMAL OR RELAXED PROBLEM} \end{array}$$

-  $h^*$  AND  $h_{\text{FF}}$  ARE PAIRWISE INCOMPARABLE

-  $h^*$  AND  $h_{\text{ADD}}$  ARE INCOMPARABLE

## - ABSTRACTION

- DROPPING SOME DISTINCTIONS BETWEEN STATES

- BUT TRYING TO PRESERVE THE TRANSITION BEHAVIOR AS MUCH AS POSSIBLE

- ABSTRACTION OF TRANSITION SYSTEM  $T$  IS DEFINED BY ABSTRACTION MAPPING  $\alpha$

- IT DEFINES WHICH STATES SHOULD BE DISTINGUISHED AND WHICH NOT

- CREATED TRANSITION SYSTEM  $T'$  IS SMALLER

- GOAL DISTANCE IN  $T'$  IS USED AS HEURISTIC

- CREATING  $T'$

- WE WANT ADMISSIBLE HEURISTIC

- SO  $h^*(\alpha(s))$  SHOULD NEVER OVERESTIMATE  $h^*(s)$

- WE HAVE TO ~~ADVANTAGE~~ ENSURE THAT ALL SOLUTIONS IN  $T$  ALSO EXIST IN  $T'$

- IF  $s$  IS GOAL STATE IN  $T$ , THEN  $\alpha(s)$  IS GOAL STATE IN  $T'$

- IF  $T$  HAS TRANSITION FROM  $s$  TO  $t$ , THEN  $T'$  HAS A TRANSITION FROM  $\alpha(s)$  TO  $\alpha(t)$



- IT ALSO MUST BE EFFICIENT COMPUTABLE

-  $\alpha(s)$

-  $\alpha^*(\alpha(s))$

- WE CAN ALSO USE MAXIMUM OF SEVERAL ADMISSIBLE HEURISTIC, WHICH CREATES NEW ADMISSIBLE HEURISTIC WHICH DOMINATES OLD ONES

- OR WE CAN ADD SEVERAL ADMISSIBLE HEURISTIC

- BUT IT CAN BECOME NON-ADMISSIBLE

- TRANSITION SYSTEM

- 5-TUPLE  $\langle S, L, T, I, G \rangle$

-  $S$  IS FINITE SET OF STATES

-  $L$  IS FINITE SET OF LABELS

-  $T \subseteq S \times L \times S$  IS TRANSITION RELATION

-  $I \subseteq S$  IS SET OF INITIAL STATES

-  $G \subseteq S$  IS SET OF GOAL STATES

-  $T$  HAS TRANSITION  $\langle s, l, s' \rangle$  IF  $\langle s, l, s' \rangle \in T$

- TRANSITION SYSTEM OF SAS<sup>+</sup> PLANNING TASK

-  $\pi = \langle V, I, O, G \rangle$  IS SAS<sup>+</sup> PLANNING TASK

- TRANSITION SYSTEM OF  $\pi$  IS

$T(H) = \langle S', L', T', I', G' \rangle$

WHERE

-  $S'$  IS SET OF STATES OVER  $V$

-  $L' = O$

-  $T' = \{ \langle s', o, s' \rangle \in S' \times L' \times S' \mid \text{APP}_O(s') = \langle \cdot \rangle \}$

-  $I' = \{ I \}$

-  $G' = \{ s' \in S' \mid s' \models G \}$

- ABSTRACTION MAPPING

-  $T = \langle S, L, T, I, G \rangle$

-  $T' = \langle S', L', T', I', G' \rangle$

-  $L = L'$

-  $\alpha \doteq S \rightarrow S'$

-  $T'$  IS ABSTRACTION OF  $T$  WITH ABSTRACTION MAPPING  $\alpha$  IF

- ~~FOR~~ FOR ALL  $s \in I$ , WE HAVE  $\alpha(s) \in I'$

- FOR ALL  $s \in G$ , WE HAVE  $\alpha(s) \in G'$

- FOR ALL  $\langle s, l, e \rangle \in T$  WE HAVE  $\langle \alpha(s), l, \alpha(e) \rangle \in T'$

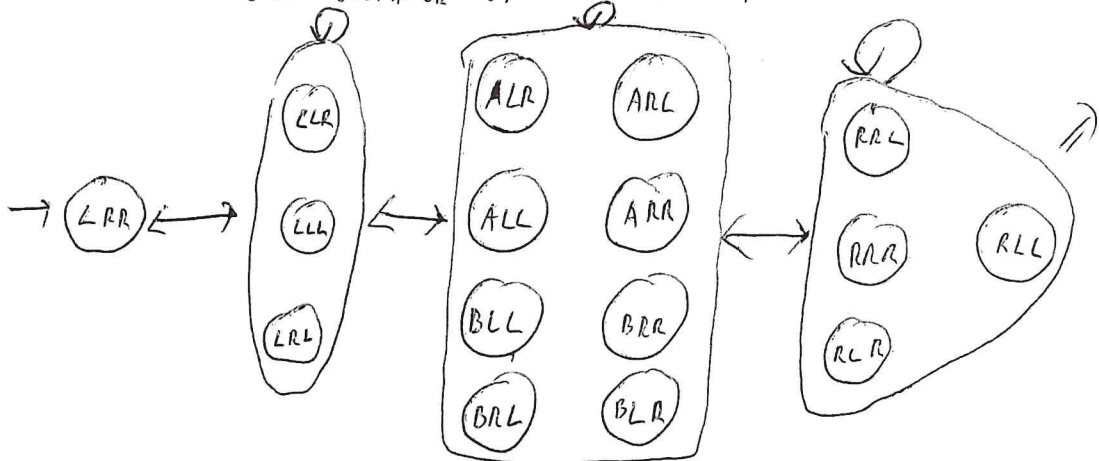
- ABSTRACTION HEURISTIC INDUCED BY  $\hat{A}$  (CONCRETE ABSTRACTION) AND  $\alpha$  (ABSTRACTION MAPPING)  $R^{A, \alpha}$

- IS HEURISTIC FUNCTION

$R^{A, \alpha} : S \rightarrow \mathbb{N}_0 \cup \{\infty\}$

WHICH MAPS EACH STATE  $s \in S$  TO  $R_A^*(\alpha(s))$

- "GOAL DISTANCE OF  $\alpha(s)$  IN  $A$ "



$R^{A, \alpha}(LRR) = 3$

- LET  $\pi$  BE SAS<sup>T</sup>

- LET  $A$  BE ABSTRACTION OF  $T(\pi)$  WITH ABSTRACTION MAPPING  $\alpha$

- THEN  $h^{A, \alpha}$  IS SAFE, GOAL AWARE, ADMISSIBLE AND CONSISTENT

- ORTHOGONAL ABSTRACTION MAPPINGS

-  $\alpha_1, \alpha_2$  BE ABSTRACTION MAPPINGS ON  $T$

-  $\alpha_1$  AND  $\alpha_2$  ARE ORTHOGONAL IF  $\#$

- FOR ALL TRANSITIONS  $\langle s, e, t \rangle$  OF  $T$

WE HAVE  $\alpha_i(s) = \alpha_i(t)$  FOR AT LEAST ONE  $i \in \{1, 2\}$

- ADDITIVITY OF ORTHOGONAL ABSTRACTION MAPPINGS

-  $h^{A_1, \alpha_1}, \dots, h^{A_n, \alpha_n}$  ARE ABSTRACTION HEURISTICS FOR SAME PLANNING TASK  $\pi$  SUCH THAT  $h^{A_i, \alpha_i}$  ARE ORTHOGONAL FOR ALL  $i \neq j$

- THEN  $\sum_{i=1}^n h^{A_i, \alpha_i}$  IS SAFE, GOAL-AWARE, ADMISSIBLE

AND CONSISTENT HEURISTIC FOR  $\pi$

- IN PRACTICE WE WANT INFORMATIVE HEURISTIC WHICH IS SMALL AT THE SAME TIME

- CONFLICTING GOAL

# - PATTERN DATABASE HEURISTICS

- ABSTRACTION HEURISTICS

- SOME ASPECTS OF TASK ARE REPRESENTED IN PERFECT PRECISION

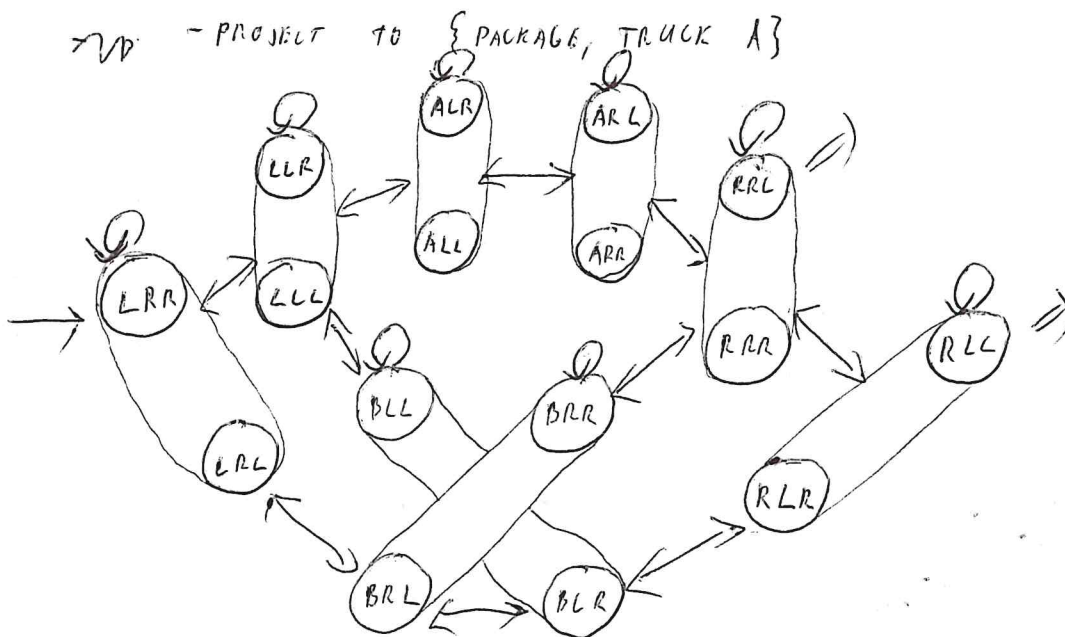
- OTHER ASPECTS ARE NOT REPRESENTED AT ALL

- EXAMPLE

- PACKAGE = {L, R, A, B}

- TRUCK A = {L, R}

- TRUCK B = {L, R}



- MERGE AND SHRINK ABSTRACTIONS

- INSTEAD OF PERFECTLY REFLECTING FEW STATE VARIABLES

- REFLECT ALL STATE VARIABLES BUT IN POTENTIALLY LOSSY WAY

- MERGE

- A AND A' CAN BE MERGED INTO PRODUCT ABSTRACTION

- SHRINK

- DUE TO MEMORY LIMITATIONS

- WE CAN USE ANY ANY ABSTRACTION ON INTERMEDIATE RESULT AND CONTINUE WITH MERGING PROCESS

## - FINITE DOMAIN REPRESENTATION

- TUPLE  $\langle V, A, I, G \rangle$

-  $V$  - FINITE SET OF STATE VARIABLES WITH FINITE DOMAINS  $\text{DOM}(V_i)$

-  $I$  - INITIAL STATE IS COMPLETE ASSIGNMENT TO  $V$

-  $G$  - GOAL IS PARTIAL ASSIGNMENT TO  $V$

-  $A$  - FINITE SET OF ACTIONS

-  $\text{PRE}(a)$ ,  $\text{EFF}(a)$

- PARTIAL ASSIGNMENTS TO  $V$

- IF COST SENSITIVE PLANNING

- EACH ACTION HAS COST  $c(a)$

## - LANDMARK

- FORMULA THAT MUST BE TRUE AT SOME POINT IN ANY PLAN

- THEY CAN BE PARTIALLY ORDERED

- SOME CAN BE DISCOVERED AUTOMATICALLY

- CURRENT APPROACHES CONSIDER ONLY LANDMARKS THAT ARE FACTS OR DISJUNCTIONS OF FACTS

## - ACTION LANDMARK

- ACTION WHICH OCCURS IN EVERY PLAN

- LANDMARKS CAN IMPLY ACTION LANDMARKS AND

- ACTION LANDMARKS CAN IMPLY LANDMARKS

- SOME ACTION LANDMARKS CAN BE DISCOVERED AUTOMATICALLY TOO

- NATURAL ORDERING

$$A \rightarrow B$$

- A IS TRUE SOME TIME BEFORE B IS TRUE

- NECESSARY ORDERING

$$A \rightarrow_n B$$

- A IS TRUE ONE STEP BEFORE B IS TRUE

- GREEDY - NECESSARY ORDERING

$$A \rightarrow_{gn} B$$

- A IS TRUE ONE STEP BEFORE B BECOMES TRUE FOR THE FIRST TIME

$$A \rightarrow_n B \Rightarrow A \rightarrow_{gn} B \Rightarrow A \rightarrow B$$

- DECIDING IF GIVEN FACT IS LANDMARK IS PSPACE - COMPLETE

- A IS LANDMARK  $\Leftrightarrow \pi'_A$  IS UNSOLVABLE

-  $\pi'_A$  IS  $\pi$  WITHOUT OPERATORS THAT ACHIEVE A

- LANDMARK DISCOVERY

- BY BACK CHAINING

- ALL GOALS ARE LANDMARKS

- IF B IS LANDMARK AND ALL ACTIONS THAT ACHIEVE B SHARE A AS PRECONDITION

- A IS LANDMARK

$$A \rightarrow_n B$$

- DISJUNCTION LANDMARK

- B IS LANDMARK

→

- ALL ACTIONS THAT ACHIEVE B HAVE A OR C  
AS PRECONDITION THEN

-  $A \vee C$  IS LANDMARK

- GENERALIZES TO ANY NUMBER OF DISJUNCTS

- DOMAIN TRANSITION GRAPHS

- DTGS

- CAN FIND LANDMARKS

- DTG OF  $N \in V$  REPRESENTS HOW THE VALUE OF  $N$   
CAN CHANGE

- DTG <sub>$N$</sub>  IS DIRECTED GRAPH WITH NODES  $D_N$  THAT  
HAS ARC  $\langle d, d' \rangle$  IFF

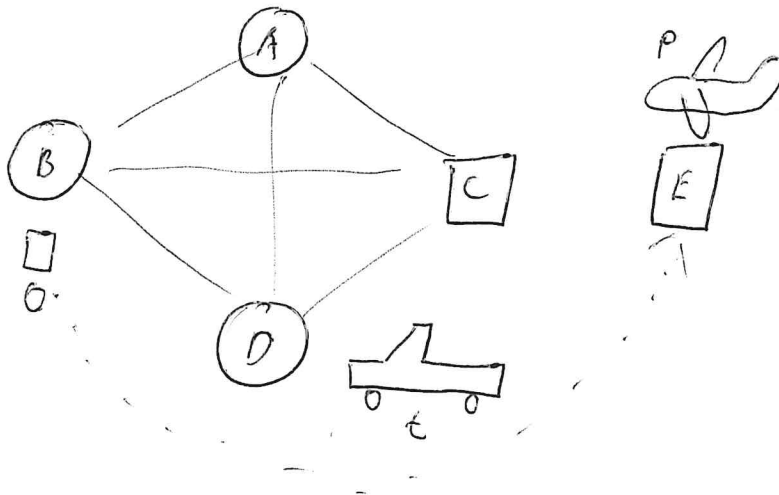
-  $d \neq d'$

- AND THERE EXISTS ACTION WITH  $N \rightarrow d'$  AS EFFECT  
AND EITHER

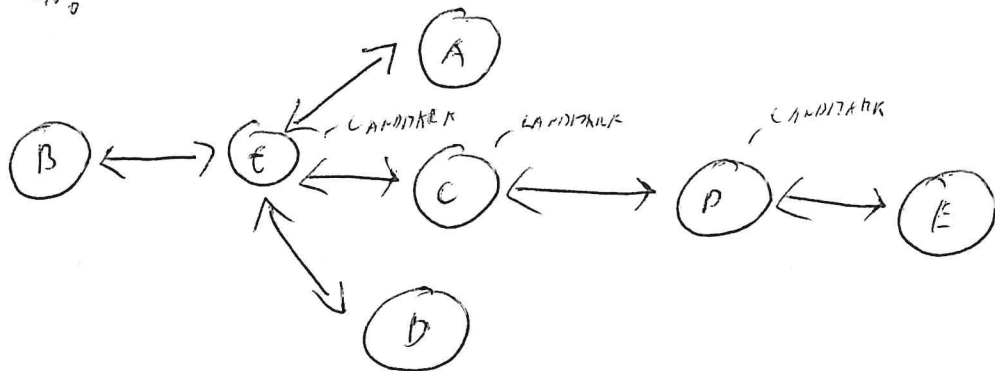
-  $N \mapsto d$  AS PRECONDITION OR

- NO PRECONDITION  $\emptyset N$

- EXAMPLE



- DTGND



- USING LANDMARKS

- DISCOVERING LANDMARK HAPPENS IN PRE PROCESSING PHASE

- LANDMARKS CAN BE USED AS SUBGOALS FOR BASE PLANNING

- PROS

- FASTER PLANNING

- HAVE TO SEARCH TO LESSER DEPTH

- CONS

- CAN LEAD TO LONGER PLANS

- NOT COMPLETE IN PRESENCE OF DEAD ENDS

- NUMBER OF LANDMARKS ~~AND~~ STILL NEEDED TO ACHIEVE GOAL  
BE USED AS HEURISTIC



## - LP - BASED HEURISTICS

### - COST PARTITIONING

- CREATE COPIES  $\pi_1, \dots, \pi_m$  OF PLANNING TASK  $\pi$

- EACH HAS ITS OWN OPERATOR COST FUNCTION

$$\text{COST}_{\pi_i}: O \rightarrow \mathbb{R}_0^+$$

- FOR ALL  $O$ : REQUIRE  $\text{COST}_{\pi_1}(O) + \dots + \text{COST}_{\pi_m}(O) \leq \text{COST}(O)$

- SUM OF SOLUTION COSTS IN COPIES IS ADMISSIBLE HEURISTIC

$$h_{\pi_1}^* + \dots + h_{\pi_m}^* \leq h_{\pi}^*$$

- WAYS TO DEFINE COST FUNCTIONS

- UNIFORM:  $\text{COST}_{\pi_i}(O) = \text{COST}(O)/m$

- ZERO-ONE: FULL OPERATOR COST IN ONE COPY, ZERO IN ALL OTHERS

### - OPTIMAL COST PARTITIONING

- USE VARIABLES FOR COST OF EACH OPERATOR IN EACH TASK COPY

- EXPRESS HEURISTIC VALUES WITH LINEAR CONSTRAINTS

- MAXIMIZE SUM OF HEURISTIC VALUES

### - LP FOR SHORTEST PATH IN STATE SPACE

MAX GOAL-DIST

SE  $\text{DISTANCE}_{S_I} = 0$  - FOR INITIAL STATE  $S_I$

$\text{DISTANCE}_{S'} \leq \text{DISTANCE}_S + \text{COST}(O)$  - FOR ALL TRANSITIONS  $S \xrightarrow{O} S'$

$\text{GOAL-DIST} \leq \text{DISTANCE}_{S^*}$  - FOR ALL GOAL STATES  $S^*$

## - OPTIMAL PARTITIONING FOR ABSTRACTIONS

$$\text{MAX } \sum_{\alpha} \text{GOAL\_DIST}^{\alpha}$$

st

FOR ALL OPERATORS  $O$

$$\sum_{\alpha} \text{COST}_O^{\alpha} \leq \text{COST}(O)$$

$$\text{COST}_O^{\alpha} \geq 0 \quad \text{FOR ALL ABSTRACTIONS } \alpha$$

FOR ALL ABSTRACTIONS  $\alpha$

$$\text{DISTANCE}_{S_I}^{\alpha} = 0 \quad \text{- FOR ABSTRACT INITIAL STATE } S_I$$

$$\text{DISTANCE}_{S'}^{\alpha} \leq \text{DISTANCE}_S^{\alpha} + \text{COST}_O^{\alpha} \quad \text{- FOR ALL TRANSITIONS } S \xrightarrow{O} S'$$

$$\text{GOAL\_DIST}^{\alpha} \leq \text{DISTANCE}_{S^*}^{\alpha} \quad \text{- FOR ALL ABSTRACT GOAL STATES } S^*$$

## - OPERATOR COUNTING

- LINEAR CONSTRAINTS WHOSE VARIABLES DENOTE NUMBER OF OCCURRENCES OF A GIVEN OPERATION

- MUST BE SATISFIED BY EVERY PLAN THAT SOLVES THE TASK

- EXAMPLES

$$Y_{O_1} + Y_{O_2} \geq 1 \quad \text{- "MUST USE } O_1 \text{ OR } O_2 \text{ AT LEAST ONCE"}$$

- ~~LP~~ LP

$$\text{MIN } \sum_O x_O \cdot \text{COST}(O)$$

$$\text{st } x_O \geq 0$$

AND SOME OPERATOR-COUNTING CONSTRAINTS

- IT IS ADMISSIBLE HEURISTIC

- SOLVABLE IN POLYNOMIAL TIME

- ADDING CONSTRAINTS MAKE HEURISTIC MORE INFORMED

## - STATE EQUATION HEURISTIC

### - SEQ

- FACTS CAN BE PRODUCED (MADE TRUE) OR CONSUMED (MADE FALSE) BY OPERATOR

- NUMBER OF PRODUCING AND CONSUMING OPERATORS MUST BALANCE OUT FOR EACH FACT

$$0 = \sum_{f \text{ PRODUCES } f} \gamma_0 - \sum_{f \text{ CONSUMES } f} \gamma_0$$

- SPECIAL CASE FOR INITIAL AND GOAL STATE

$$G(f) - S(f) = \sum_{f \text{ PRODUCES } f} \gamma_0 - \sum_{f \text{ CONSUMES } f} \gamma_0$$

## - POTENTIAL HEURISTICS

- HEURISTIC DESIGN AS OPTIMIZATION PROBLEM

- DEFINE SIMPLE NUMERICAL STATE FEATURES  $f_1, \dots, f_m$

- CONSIDER HEURISTICS THAT ARE LINEAR COMBINATIONS OF FEATURES

$$h(s) = w_1 f_1(s) + \dots + w_m f_m(s)$$

-  $w_m$  IS POTENTIAL

- GOAL IS TO FIND POTENTIALS ( $w_i$ ) FOR WHICH  $h$  IS ADMISSIBLE AND WELL INFORMED

- SUCH HEURISTIC IS VERY FAST TO COMPUTE IF ~~HEURISTICS~~ FEATURES ARE

## - FEATURES

- FEATURE IS NUMERICAL FUNCTION DEFINED ON STATES OF TASK

$$f: S \rightarrow \mathbb{R}$$

## - ATOMIC FEATURE

$$- f_{X=x}(s) = \begin{cases} 1 & \text{IF VARIABLE } X \text{ HAS VALUE } x \text{ IN STATE } s \\ 0 & \text{OTHERWISE} \end{cases}$$

- EXAMPLE  $R(s) = 3 f_{X=a} + \frac{1}{2} f_{X=b} - 2 f_{X=c} + \frac{5}{2} f_{X=d}$

- HOW SET WEIGHT

## - REQUIREMENTS

- WE WANT HEURISTIC WHICH IS

- ADMISSIBLE

- CONSISTENT

- WELL INFORMED

## - LINEAR PROGRAMMING

- CONSTRAINTS

- GOAL-AWARENESS ( $R(s) = 0$  FOR GOAL STATE  $s$ )

$$\sum_{\text{GOAL FACTS } f} w_f = 0$$

- CONSISTENCY

$$\sum_{\substack{f \text{ CONSUMED} \\ B \geq 0}} w_f - \sum_{\substack{f \text{ PRODUCED} \\ B \geq 0}} w_f \leq \text{COST}(o) \text{ FOR ALL OPERATORS } o$$

- IF GOAL-AWARE AND CONSISTENT  $\Rightarrow$  ADMISSIBLE AND CONSISTENT

- WELL INFORMED

- ENCODE QUALITY METRIC INTO OBJECTIVE FUNCTION  
AND MAXIMIZE IT

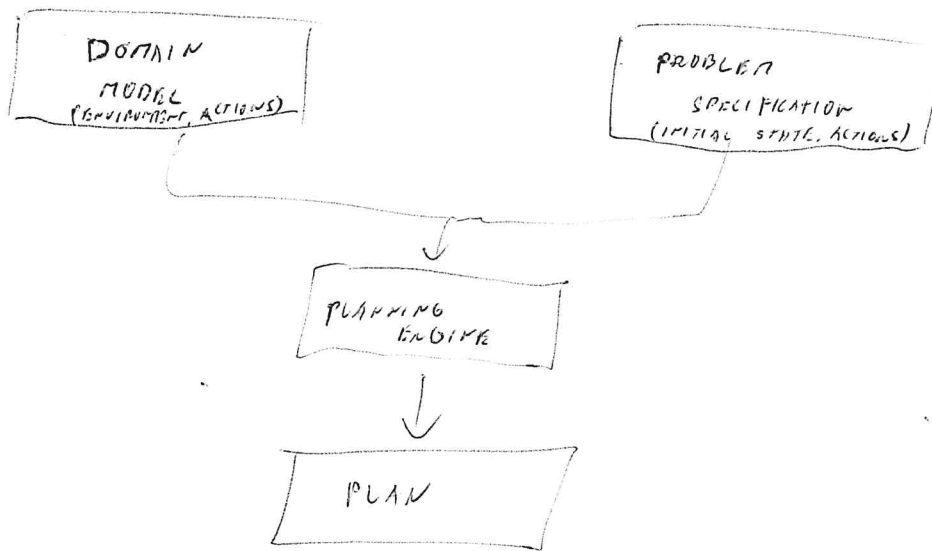
- EXAMPLES

- MAXIMIZE HEURISTIC VALUE OF GIVEN STATE (INITIAL  
FOR EXAMPLE)

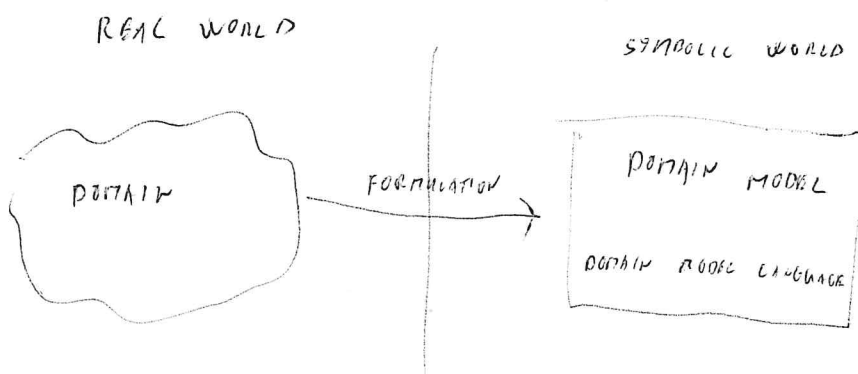
- MAXIMIZE AVERAGE HEURISTIC VALUE OF ALL (OR SOME)  
STATES

- MINIMIZE ESTIMATED SEARCH EFFORT

- DOMAIN INDEPENDENT PLANNING CONCEPT



91



## - MODELING LANGUAGES

- PDDL

- MOST WIDESPREAD

- NDDL

- APTL

- RDDL

## - DOMAIN CONTROL KNOWLEDGE

- DCK

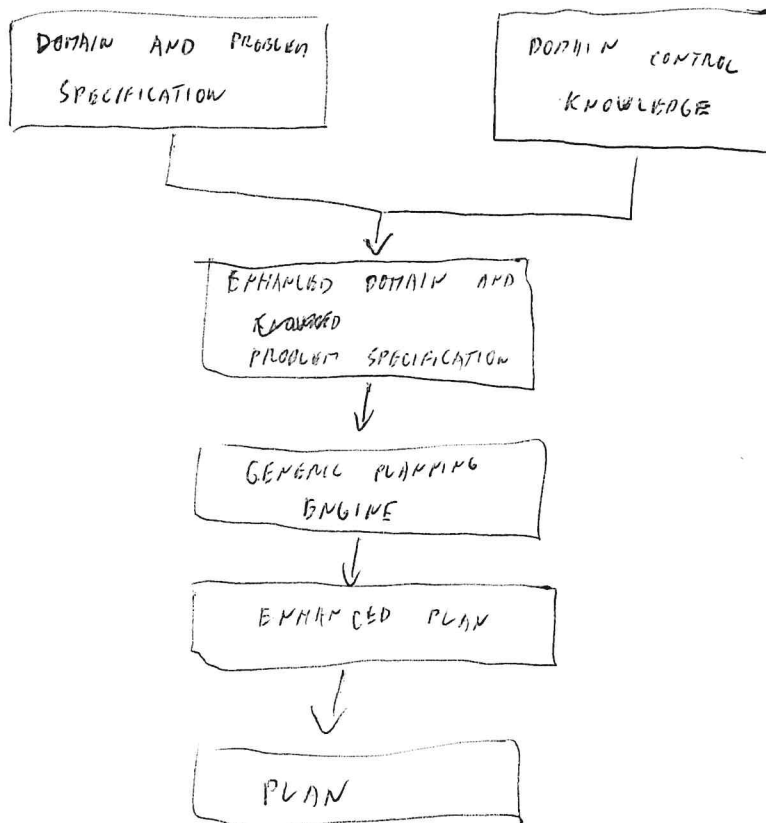
- CAPTURES USEFUL DOMAIN-SPECIFIC INFORMATION

- PROVIDES GUIDANCE FOR PLANNING ENGINES

- CAN BE

- PLANNER SPECIFIC

- PLANNER INDEPENDENT



## - TEMPORAL PLANNING

- PLANNING IN SITUATIONS WHERE ACTIONS

- HAVE NON-ZERO DURATION

- MAY OVERLAP IN TIME

- WE NEED REPRESENTATION OF ~~REAL~~ TIME

- STATE VARIABLE

- PARTIALLY SPECIFIED FUNCTION TELLING WHAT IS TRUE AT SOME TIME  $t$

- TEMPORAL ASSERTIONS

- EVENT

-  $x @ t : (v_1, v_2)$

- AT TIME  $t$ ,  $x$  CHANGES FROM  $v_1$  TO  $v_2 \neq v_1$

- PERSISTENCE CONDITION

-  $x @ [t_1, t_2) : v$

-  $x = v$  THROUGHOUT INTERVAL  $[t_1, t_2)$

- WHERE

$t_i$  ARE ~~CONSTANTS~~ <sup>CONSTANTS</sup> OR TEMPORAL VARIABLES

$v_i$  ARE CONSTANTS OR OBJECT VARIABLES

# - PLANNING IN MULTIAGENT SYSTEMS

## - MULTIAGENT SYSTEM

### - CONSISTS OF AGENTS

- BOUNDED ENTITY

- INTERACTING WITH OTHERS

### - AGENT

#### - REACTIVE AGENT

- REACT TO IMMEDIATE STIMULI

#### - PROACTIVE AGENT

- PROACTIVELY TRIES TO SATISFY ITS NEEDS AND GOALS

#### - PLANNING AGENT

- PLANS ITS ACTIONS IN ADVANCE

### - HOW SEES OTHER AGENTS

#### - ADVERSARY

- GAME THEORY

- PART OF ENVIRONMENT

- PONDOP

- COOPERATION

## - MULTIAGENT PLANNING PROBLEM

$$- M = \langle A, \{\pi_i\}_{i=1}^n \rangle$$

- A ... SET OF AGENTS

-  $\{\pi_i\}_{i=1}^n$  ... SET OF  $n$  AGENT PLANNING PROBLEMS

$$- \pi_i = \langle V_i, O_i, S_{\emptyset}^i, S_{\neq} \rangle$$

-  $V_i = V^{\text{PUB}} \cup V_i^{\text{PRIV}}$  - SET OF PUBLIC AND PRIVATE VARIABLES

-  $O_i = O_i^{\text{PUB}} \cup O_i^{\text{PRIV}}$  - SET OF PUBLIC AND PRIVATE OPERATIONS



-  $S_i^j$  - INITIAL STATE OF AGENT  $i$   
(FULL VALUATION OF  $V_i$ )

-  $S_*$  - PUBLIC GOAL STATE (PARTIAL VALUATION  
OVER  $V_{pub}$ )

- SOLUTION

$$\pi_1 = (o_1, \epsilon, \epsilon, \dots, o_2, o_3)$$

$\vdots$

$$\pi_n = (\epsilon, o_4, \epsilon, \dots, o_5, \epsilon)$$

-  $\epsilon$  IS NO-OPTION

- ACTIONS IN EACH TIME-STEP MUST NOT BE MUTUALLY  
EXCLUSIVE

- GLOBAL PROBLEM

$$\pi^G = \langle V = \bigcup_{i=1}^n V_i, O = \bigcup_{i=1}^n O_i, S_0, S_* \rangle$$

- WHY CAN'T WE SOLVE GLOBAL PROBLEM CENTRALLY

- COST OF COMMUNICATION

- COST OF FORMALIZATION

- SPEED IMPROVEMENT

- PRIVACY

- WEAK PRIVACY

- DO NOT SHARE PRIVATE INFORMATION

- BUT IT CAN BE INFERRRED FROM PUBLIC INFORMATION AND  
EXECUTION OF ALGORITHM

- STRONG PRIVACY

- NO PRIVATE INFORMATION CAN BE INFERRRED FROM PUBLIC  
INFORMATION AND RUN OF ALGORITHM

- HARD TO ACHIEVE

- MAP PLANNING PARADIGMS

- PLAY-BASED COORDINATION

- STATE-BASED COORDINATION

- PLAN-BASED COORDINATION

- LET  $\pi^{\Delta} = (o_1^{\Delta}, \dots, o_n^{\Delta})$  BE A PUBLIC PLAN

- WE SAY THAT  $\pi^{\Delta}$  IS EXTENSIBLE BY AGENT  $i$

- IF BY INSERTING SOME  $o_{i,1}, \dots, o_{i,j} \in O_i^{\text{priv}}$  INTO RESPECTIVE  $\pi = (o_1, \dots, o_n)$  WE OBTAIN  $\pi'$  WHICH IS A SOLUTION TO THE AGENT PROJECTION  $\pi^{\Delta, i}$

- LET  $\pi^{\Delta}$  BE A PUBLIC PLAN. IF  $\pi^{\Delta}$  IS EXTENSIBLE BY ALL AGENTS IN  $A$  THEN  $\pi^{\Delta}$  IS A PUBLIC PROJECTION OF SOLUTION  $\pi$  OF  $M$

- GENERAL IDEA

- FIND PUBLIC PLAN EXTENSIBLE BY ALL AGENTS

- FIND RESPECTIVE EXTENDING PLAN FOR EACH AGENT

- PLANNING FIRST

- DISTRIBUTED CSP TO FIND EXTENSIBLE PUBLIC PLAN

- ONE VARIABLE PER AGENT

- VALUES

- ALL POSSIBLE SEQUENCES OF PUBLIC ACTIONS FROM  $O^{\Delta}$  AND PUNCHEDERS FOR PRIVATE ACTIONS

- CONSTRAINTS

- BINARY, SEQUENCES MUST MATCH ON PUBLIC ACTIONS

- UNITY; THE VALUE (PUBLIC PLAN) MUST BE EXTENSIBLE BY RESPECTIVE

- PSM

- EACH AGENT  $i$

- GENERATE PLAN FOR  $\pi^{D_i}$

- CREATE PUBLIC PROJECTIONS

- SHARE AND FIND INTERSECTION

- IF INTERSECTION NONEMPTY, SOLUTION FOUND

- ELSE ITERATE

- STATE BASED COORDINATION

- MULTI AGENT FORWARD SEARCH

- SEMANTIC SEARCH SPACES

- EACH AGENT SEARCH USING OPERATORS IN  $O_i$

- IF STATE  $s$  WAS REACHED BY PUBLIC OPERATOR

$OE_i^{pub}$ , SEND  $s$  TO ALL OTHER AGENTS

- ENCRYPT PRIVATE PARTS OF STATE

- MAD- $A^*$

- HOW ENSURE GLOBAL OPTIMALITY?

- BETTER STATE MIGHT BE IN OTHER

AGENTS' OPEN LIST

- WE NEED TO CHECK GLOBAL STATE OF  
DISTRIBUTED SYSTEM

## - HEURISTIC IN MULTI AGENT PLANNING

### - PROJECTED HEURISTIC

- COMPUTED ON THE AGENT PROJECTED PROBLEM BY ESTIMATING ONLY AGENTS PROBLEM SOLUTION COST

### - MAINTAINS ADMISSIBILITY

- AGENTS PROJECTION IS ABSTRACTION

- TAKING MAX  $(h_i(s), h_j(s))$  IS ADMISSIBLE

### - STRONG PRIVACY

- NO INFORMATION EXCHANGES

- UNLESS HEURISTIC ITSELF LEAKS INFO

### - DISTRIBUTED HEURISTIC

- COMPUTED IN DISTRIBUTED WAY

- ESTIMATES GLOBAL PROBLEM SOLUTION COST

- BETTER INFORMED

→ ADDRESS

### - POTENTIAL HEURISTIC

$$- h_{POT}(s) = \sum_{N \in V} POT(\langle V, S[V] \rangle)$$

### - IN MAP

$$h_{POT}^{PUB}(s) = \sum_{N \in V^{PUB}} POT(\langle V, S[V] \rangle)$$

$$h_{POT}^{PRIV_i}(s) = \sum_{N \in V_i^{PRIV}} POT(\langle V, S[V] \rangle)$$

$$h_{POT}^G(s) = h_{POT}^{PUB}(s) + \sum_{i \in A} h_{POT}^{PRIV_i}(s)$$

- NO ADDITIONAL COMMUNICATION

- PRIVATE PARTS  $R_{POT}^{PRIV}(s)$  OF OTHER AGENTS  $i \in A \setminus \{j\}$  ARE NOT CHANGED BY AGENT  $j$

## - PROBABILISTIC PLANNING

- ACTIONS TAKE SOME TIME

- ACTIONS HAVE NON-DETERMINISTIC RESULTS

- ENVIRONMENT MAY CHANGE

- AGENT DOESN'T KNOW WHOLE SITUATION

- NOT-PRECISE SENSORS

## - CLASSICAL PLANNING

$$\langle S, s_0, S_G, A, f, c \rangle$$

-  $f$  ... TRANSITION FUNCTION  $f: S \times A \rightarrow S$

## - PROBABILISTIC PLANNING

- PROBABILISTIC TRANSITION FUNCTION

$$T: S \times A \times S \rightarrow [0, 1]$$

$$\sum_{s' \in S} T(s, a, s') = 1$$

- SOLUTION IN CLASSICAL PLANNING IS SET OF ACTIONS

- SOLUTION IN PROBABILISTIC PLANNING IS POLICY

$$- \pi: H \times A \rightarrow [0, 1]$$

$$h = s_1 a_1 s_2 a_2 \dots s_k$$

## - EVALUATION

- COST ASSIGNED TO  $(s, a, s')$

- "REWARD"

- EXECUTING AN POLICY YIELDS A SEQUENCE OF REWARDS

- POLICY VALUE

$$- v(R_1, R_2, \dots) = R_1 + \gamma R_2 + \gamma^2 R_3 + \dots$$

$$- v(\pi(s_0)) = E [v(R_1, \dots)]$$

- THERE EXIST POLICY THAT IS OPTIMAL AT EVERY TIME STEP

- WE SEARCH FOR

$$\pi^* \text{ s.t. } v(\pi^*) \geq v(\pi) \text{ FOR ALL OTHER } \pi$$

## - MARKOV DECISION PROCESS

-  $\langle S, A, D, T, R \rangle$

- S - FINITE SET OF STATES

- A - FINITE SET OF ACTIONS

- D - HORIZONT, FINITE/INFINITE SET OF TIMESTEPS

- T - TRANSITION FUNCTION

$$T: S \times A \times S \rightarrow [0, 1]; \sum_{s' \in S} T(s, a, s') = 1$$

- R - REWARD FUNCTION

$$R: S \times A \times S \rightarrow \mathbb{R}$$

- TYPICALLY BOUNDED

- MDP POLICY

- HISTORY DEPENDENT POLICY

$$-\pi : H \times A \rightarrow [0, 1]$$

$$\sum_{a \in A} \pi(R, a) = 1$$

- BUT FOR SIMPLE CASES WE DON'T NEED HISTORY AND RANDOMIZATION

- MARKOV ASSUMPTION

- POLICY IS ASSIGNMENT OF ACTION IN EACH STATE AND TIME

$$-\pi : S \rightarrow A$$

- STATIONARY POLICY

- POLICY IS SAME EVERY TIME  $S$  IS VISITED

- POSITIONAL POLICY

- DETERMINISTIC AND STATIONARY POLICY

- FF - REPLAN

- 1. DETERMINIZE THE DOMAIN  
(REMOVE STOCHASTICITY)

- 2. SYNTHESIZE PLAN

- 3. EXECUTE PLAN

- 4. IF UNEXPECTED STATE OCCUR, REPLAN

- WAYS OF DETERMINIZATION

- KEEP ONLY MOST PROBABLE OUTCOME OF ACTION

- KEEP ALL OUTCOMES BUT GENERATE SEPARATE ACTION FOR EACH

- SIMPLE, NOT SOUND, NOT OPTIMAL, BUT ENOUGH FOR SIMPLE DOMAINS

## - ROBUST -FF

- GENERALIZES FF-REPLAN

- 1. DETERMINIZE PROBLEM

- 2. USE CLASSICAL PLANNER TO FIND PARTIAL PLANS

- 3. AGGREGATE THESE PLANS INTO PARTIAL POLICY

- 4. CONTINUE UNTIL PROBABILITY OF REPLANNING IS BELOW GIVEN THRESHOLD

## - SEVERAL OPTIONS

- DETERMINIZATION

- SELECTING GOAL

- CALCULATING PROBABILITY OF REACHING TERMINAL STATE

- IT IS SOLVED ONLY WITH SELECTED OPTIONS

- NOT OPTIMAL IN GENERAL

## - FF- HINDSIGHT

- APPROXIMATE VALUE OF STATE

- SAMPLE SET OF DETERMINIZED PROBLEMS ORIGINATING FROM STATE

- SOLVE THESE PROBLEMS AND COMBINE THEIR VALUE

- OPTIMAL VALUE

$$V^*(s, T) = \max_{\pi} E[R(s, F, \pi)]$$

-  $s$  - STATE

-  $T$  - HORIZONT

-  $\pi$  - POLICY

-  $F$  - RANDOM VARIABLE

-  $R$  - REWARD FUNCTION



- HOP VALUE APPROXIMATION

$$V^*(s, T) = E \left[ \max_{\pi} R(s, \pi) \right]$$

- VALUE OF POLICY

$$V_{\pi}^{\lambda}(s) = E \left[ \sum_{t=0}^{\lambda} \gamma^t \cdot R(s_t, a_t, s_{t+1}) \mid s_0 = s, a_t = \pi(s_t) \right]$$

- OPTIMAL POLICY

$$\pi^{* \lambda}(s) = \underset{\pi}{\text{ARGMAX}} V_{\pi}^{\lambda}(s)$$

- FOR LARGE OR INFINITE  $\lambda$  WE CAN APPROXIMATE VALUE BY

DYNAMIC PROGRAMMING

$$- V_{\pi}^0(s) = 0$$

$$- V_{\pi}^{\lambda}(s) = \sum_{s' \in S} T(s, a, s') [R(s, a, s') + \gamma V_{\pi}^{\lambda-1}(s')] \quad a = \pi(s)$$

- VALUE ITERATION

$$- V^0(s) = 0 \quad \forall s \in S$$

$$- V^{\lambda}(s) = \max_{a \in A} \sum_{s' \in S} T(s, a, s') [R(s, a, s') + \gamma V^{\lambda-1}(s')]$$

Q-FUNCTION  $Q(s, a)$

- FOR  $\lambda \rightarrow \infty$  VALUES CONVERGE TO OPTIMUM  $V^{\lambda} \rightarrow V^*$

- IT CONVERGES

- WE CAN MEASURE RESIDUAL  $r$  AND STOP IF IT IS SMALL ENOUGH

$$(r \leq \epsilon (1 - \gamma) / \gamma)$$

$$- r = \max_{s \in S} |V_{i+1}(s) - V_i(s)|$$

- CONVERGENCE DEPENDS ON  $\gamma$

- VALUE ITERATION CALCULATES ONLY VALUES OF STATES

- OPTIMAL POLICY IS EXTRACTED BY GREEDY APPROACH

$$- \pi^{\lambda}(s) = \text{ARG MAX}_{\alpha \in A} \sum_{s' \in S} T^{\lambda}(s, \alpha, s') [R^{\lambda}(s, \alpha, s') + \gamma V^{\lambda}(s')]$$

- POLICY ITERATION

- START WITH ARBITRARY POLICY

- TWO STEPS

- POLICY EVALUATION

- RECALCULATES VALUE OF STATES GIVEN CURRENT POLICY  $\pi^{\lambda}$

- POLICY IMPROVEMENT

- CALCULATE NEW MAXIMUM EXPECTED UTILITY POLICY  $\pi^{\lambda+1}$

- ITERATE UNTILL VALUE CHANGES

- ASYNCHRONOUS VALUE ITERATION

- YOU DON'T HAVE TO UPDATE ALL STATES AT ONCE

- YOU CAN UPDATE ONE STATE AT A TIME

- THIS LOWERS THE MEMORY CONSUMPTION

- WE CAN INITIALIZE VALUES OF STATES MORE INTELLIGENTLY THAN BY 0

- WE CAN RUN FF-REPLAN ON THEM FOR EXAMPLE

# - MONTE CARLO METHODS

- REPEATED SAMPLING TO DETERMINE THE PROPERTIES OF SOME PHENOMENON

## - MONTE CARLO PLANNING

- COMPUTE GOOD POLICY FOR MDP BY INTERACTING WITH MDP SIMULATOR

## - UCB 1

- SELECTION ACTION FORMULA

- SELECT ACTION  $a_i$  WHICH MAXIMIZES

$$g_i^{u+c} \sqrt{\frac{\ln n}{n_i}}$$

AND UPDATE  $g_i$

-  $n$  ... NUMBER OF TIMES THE STATE IS VISITED

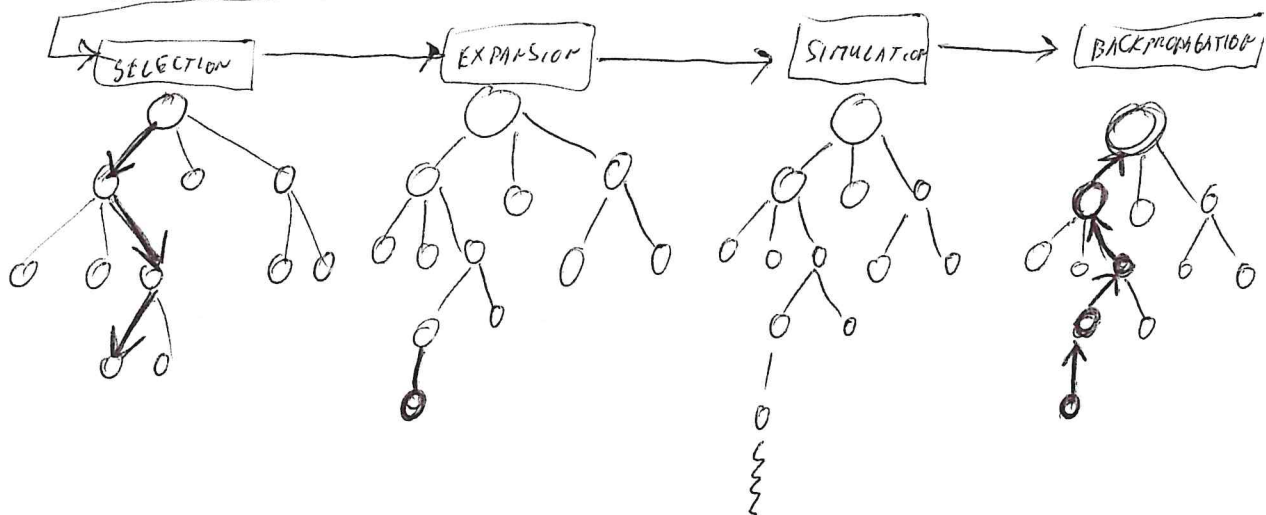
-  $n_i$  ... NUMBER OF TIMES THE ACTION IS VISITED

-  $g_i$  ... AVERAGE REWARD FROM PREVIOUS PLAYS

-  $c$  ... EXPLORATION FACTOR, ENSURES THAT RARELY EVALUATED ACTIONS ARE EVALUATED

## - MONTE CARLO TREE SEARCH

- UCT



- VANILA UCT DOESN'T WORK WELL IN PRACTICE

- HUGE BRANCHING FACTOR

- LONG HORIZONT

- DIFFICULT TO FIND CORRECT PLAY BY RANDOM POLICIES

- PROST

- SEARCH DEPTH LIMITATION

- PRUNING OUT UNREASONABLE ACTIONS

- HEURISTIC VALUE INITIALIZATION

- ONLINE PLANNING

- SELECTING BEST ACTION IN CURRENT SITUATION IN LIMITED TIME

- SIMPLE REGRET

- WE DON'T WANT TO REGRET NOT SELECTING DIFFERENT ACTION IN CURRENT STATE

- BUT UCB 1 OPTIMIZES CUMULATIVE REGRET

- SELECTING BEST ARM OVER ALL ATTEMPTS

- BRUE

- TWO CONFLICTING TASKS

- SELECT THE BEST ACTION IN  $s$  (REACHING  $s$ )

- EXPLORING AND FINDING THE BEST CONTINUATION AFTER  $s$  IS REACHED

- TO SATISFY SECOND TASK

- WE HAVE TO SELECT THE BEST ACTION SUFFICIENTLY OFTEN

- BUT TO DO THAT WE NEED TO KNOW OPTIMAL CONTINUATION

- BRUE USES TWO ACTION SELECTION METHODS

- SELECTION PHASE

- ACTION IS CHOSED UNIFORMLY

- IF UPDATE PHASE

- ACTION IS SELECTED USING GREEDY STRATEGY

- TRIAL-BASED HEURISTIC TREE SEARCH

- COMMON FRAMEWORK WITH FIVE PARTS

- HEURISTIC FUNCTION

- BACKUP FUNCTION

- ACTION SELECTION

- OUTCOME SELECTION

- TRIAL LENGTH

- MAINTAINS TREE OF ALTERNATING DECISION AND CHANCE NODES

- SELECTION PHASE

- ALTERNATING VISIT\_DECISION\_MODE AND VISIT\_CHANCE\_MODE

- SELECTING ACTION BY SELECT\_ACTION AND SELECT\_OUTCOME

- TRAVERSING TREE DOWN

- EXPANSION PHASE

- UNVISITED NODE IS ENCOUNTERED

- ADD CHILD FOR EACH ACTION

- USE HEURISTIC TO INITIALIZE THE ESTIMATE

- ALLOWS SELECTION PHASE FOR NEW NODES

- SELECTION AND EXPANSION ALTERNATE UNTIL TRIAL LENGTH

- BACKUP PHASE

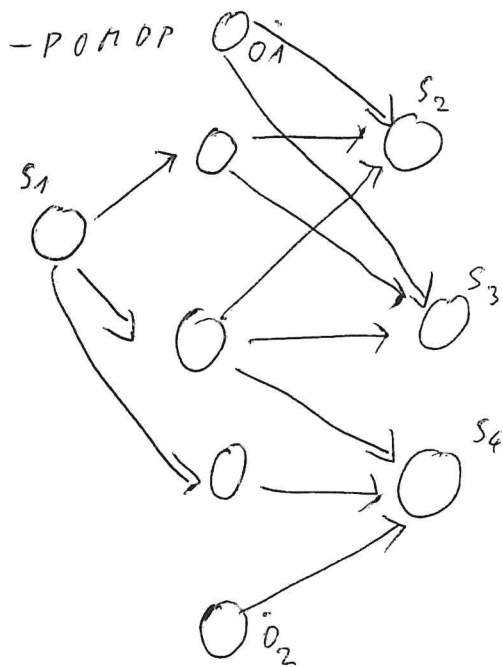
- BACKUP\_DECISION\_MODE

- BACKUP\_CHANCE\_MODE

- ALL SELECTED NODES ARE UPDATED IN REVERSE ORDER

- TRIAL ENDS WHEN BACKUP IS CALLED ON ROOT NODE

- PARTIALLY OBSERVABLE MARKOV DECISION PROCESSES



- MODEL FOR SCENARIOS WITH UNCERTAIN OBSERVATIONS

$$- \langle S, A, D, O, b_0, T, \Omega, R, \gamma \rangle$$

-  $S$  ... STATES (FINITE SET)

-  $A$  ... FINITE SET OF ACTIONS

-  $D$  ... TIME STEPS

-  $O$  ... FINITE SET OF POSSIBLE OBSERVATIONS

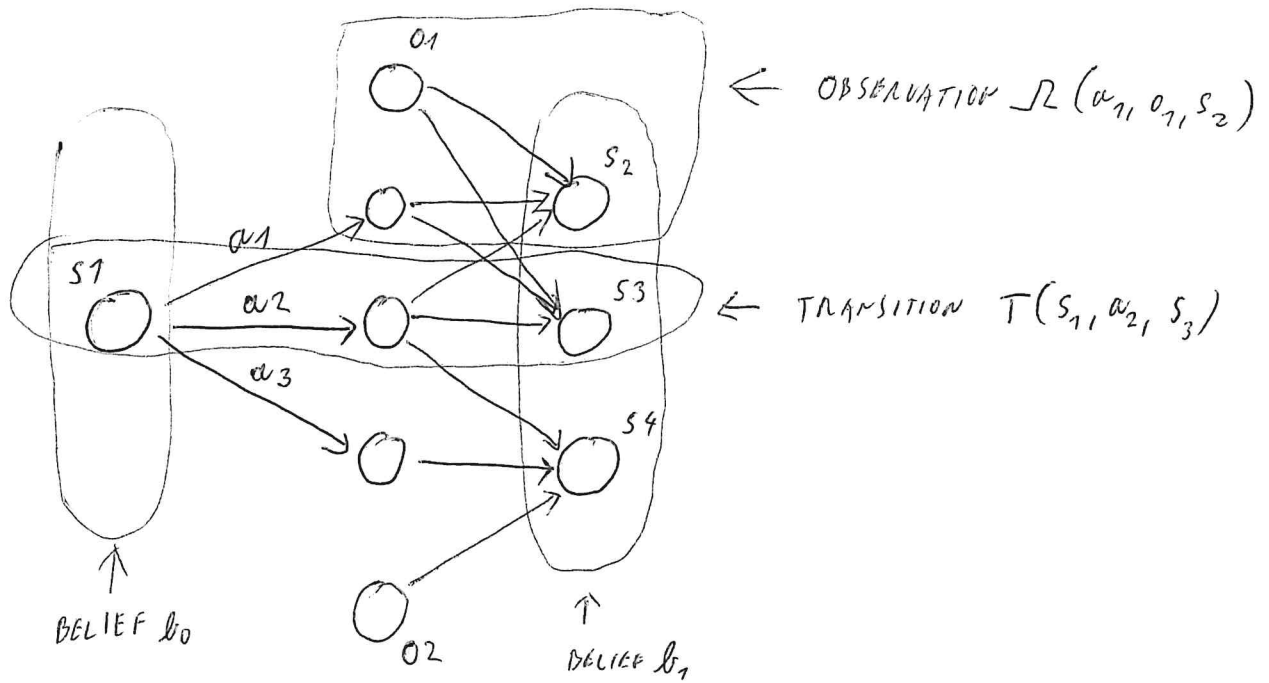
-  $b_0$  ... INITIAL BELIEF FUNCTION  $b_0: S \rightarrow [0, 1]$

-  $T$  ... TRANSITION FUNCTION  $T: S \times A \times S \rightarrow [0, 1]$

-  $\Omega$  ... OBSERVATION PROBABILITY:  $\Omega: A \times O \times S \rightarrow [0, 1]$

-  $R$  ... REWARD FUNCTION  $R: S \times A \rightarrow \mathbb{R}$

-  $\gamma$  ... DISCOUNT FACTOR  $0 \leq \gamma \leq 1$



- BELIEF IS PROBABILITY DISTRIBUTION OVER STATES

- BELIEFS ARE UNIQUELY IDENTIFIED BY HISTORY

-  $b_1$  - PROBABILITY DISTRIBUTION OVER STATES AFTER PLAYING ONE ACTION

$$- b_t \leftarrow \Pr(s_t | b_0, a_0, o_1, \dots, o_{t-1}, a_{t-1}, o_t)$$

- BY DYNAMIC PROGRAMMING

$$b_t(s') = \gamma \sum_{\omega, o, s} R(\omega, o, s') \cdot \sum_{s \in S} T(s, a, s') b_{t-1}(s)$$

-  $o \dots$  LAST OBSERVATION

-  $a \dots$  LAST ACTION

-  $\gamma \dots$  NORMALIZING CONSTANT

- BELIEFS DETERMINE NEW VALUES

$$V(b) = \max_{a \in A} [R(b, a) + \gamma \sum_{b' \in B} T(b, a, b') V(b')] ]$$

- THIS WAY WE TRANSFORMED POMDP TO MDP

- IN THEORY WE CAN USE ANY MDP ALGORITHM ~~AND~~

- BUT  $B$  IS INFINITE

# - SOLVING CONTINUOUS STATE MDPs

- VALUES CAN BE COMPACTLY REPRESENTED AS FINITE SET OF  $\alpha$  VECTORS

$$V = \{ \alpha_0, \dots, \alpha_m \}$$

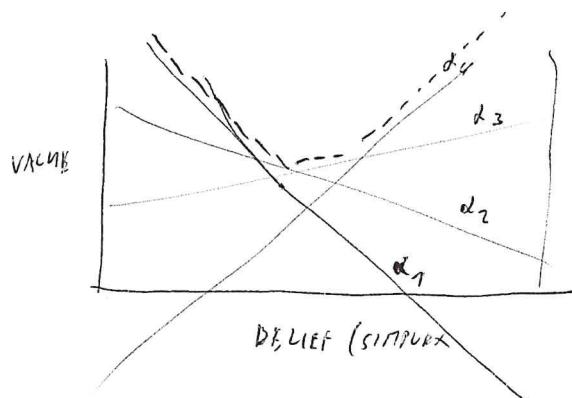
- $\alpha$  VECTOR IS  $|S|$  DIMENSIONAL HYPER-PLANE

- LINEAR FUNCTION REPRESENTING UTILITY VALUES AFTER SELECTING SOME FIXED ACTION

- DEFINES VALUE FUNCTION OVER BOUNDED REGION OF BELIEF

$$- V(B) = \max_{\alpha \in V} \sum_{s \in S} \alpha(s) b(s)$$

- $V$  IS PIECE-WISE LINEAR CONVEX FUNCTION



- MODIFICATION OF ITERATION ALGORITHM TO WORK WITH  $\alpha$

$$V^t(B) = \max_{a \in A} \left[ \sum_{s \in S} R(s, a) b(s) + \gamma \sum_{o \in O} \max_{\alpha' \in V^{t-1}} \sum_{s \in S} \sum_{s' \in S} T(s, a, s') \Omega(o, s, a) \cdot \alpha'(s') b(s) \right]$$

- FORMULA TO COMPUTE  $\alpha$ -VECTORS

$$\alpha^{a, *}(s) = R(s, a)$$

$$\alpha_i^{a, 0}(s) = \gamma \sum_{s' \in S} T(s, a, s') \Omega(o, s, a) \alpha_i'(s') \quad \forall \alpha_i' \in V'$$

$$V^a = \alpha^{a, *} \oplus \alpha^{a, 0_1} \oplus \alpha^{a, 0_2} \oplus \dots$$

$$V = \bigcup_{a \in A} V^a$$



- BUT THIS HAS SEVERAL DISADVANTAGES

- COMPLEXITY

- EXPONENTIAL IN  $|V|^{101}$

- SO USEFUL FOR SMALL DOMAINS

- POINT BASED VALUE ITERATION

- USE ONLY LIMITED SET OF BELIEFS

$$B = \{b_0, \dots, b_q\}$$

- KEEP ONLY SINGLE ALPHA VECTOR FOR ONE BELIEF POINT

- AT TIME ALGORITHM ALTERNATES

- BELIEF POINT VALUE UPDATE

- BELIEF POINT SET EXPANSION

- BELIEF VALUE UPDATE

$$V_b^{\alpha} = \alpha^{a_i, x} + \gamma \sum_{o \in O} \text{ARG MAX}_{\alpha \in \alpha_i^{a_i, o}} (\alpha_i; b)$$

$$V \leftarrow \text{ARG MAX}_{V_b^{\alpha}, b \in B} V_b^{\alpha}, \forall b \in B$$

- THIS REMOVES EXPONENTIAL COMPLEXITY

- VALUE ITERATION ENDS AFTER  $k$  ITERATIONS

- BELIEF POINT SET EXPANSION

- SAMPLING NEW BELIEFS FROM EXISTING BELIEFS

- TRYING TO UNIFORMLY COVER REACHABLE BELIEF SPACE

