



BEYOND “KNOWING THAT” (IV)

LOGICS OF KNOWING HOW

Yanjing Wang

Department of Philosophy, Peking University

NASSLLI 2018, CMU

www.wangyanjing.com

Background

Plan-based knowing how

Strategy-based knowing how

Further directions

BACKGROUND

- In epistemology:
 - Can knowledge-how be reduced to knowledge-that?
 - Anti-intellectualism: No, Knowledge-how is similar to ability (e.g., [Ryle 49])
 - Intellectualism: Yes, it is reducible based on linguistic formulation (e.g., [Stanley & Williamson 2001])
- In imperfect information games
 - Can a group of agents know how to win the game?
- In automated planning under uncertainty
 - Can an autonomous agent know how to achieve some goal?

Is knowledge-how just some ability? Can you say:

- ?I know how to digest.
- ?I know how to lift a 5kg bag.
- ?An infant knows how to ask for food.
- ?A dog knows how to catch a frisbee.
- ?A computer knows how to translate this sentence.
- ?A monkey played Chopin by luck. Does it know how?
- ?What about a well-trained piano monkey?
- ?A skier escaped the avalanche, he knows how to do it.
- ?A broken-arm pianist knows how to play piano.

Is the experience necessary in knowledge-how? Can you say:

- ?The trainer of an Olympic gym champion knows how to do the champion moves.
- ?You know the rules of Chess thus you know how to play.
- ?You know how to go to the central station even when you have never been there.
- ?A pilot knows how to fly a plane even if he was only trained in (extremely realistic) simulator.
- ?You can cook the right dish using wrong recipe and wrong ingredients (which happen to cancel each other's effects).
- ...

Clarifications:

- We do not focus on the philosophical debate between intellectualism (e.g., Stanley & Williamson 2001) and anti-intellectualism (e.g., Ryle 49). See the collection of 200+ papers on the topic at philpapers.org.
- We focus on *goal-directed* “knowing how”: knowing how to realize a goal, e.g., I know how to go to Beijing; I know how to know the answer; I know how to prove the theorem.
- We do not study “knowing how” in the following senses: I know how turtles reproduce; I know how happy she is; I know how to speak Chinese; I know how to behave at the dinner table....

Vendler's 4 categories of verbs denoting:
states, activities, accomplishments and achievements.

Dowty gives the following examples:

States	Activities	Accomplishments	Achievements
know	run	build	recognize
believe	walk	make a chair	find
have	swim	recover from illness	die

Activity directed, rule directed, goal directed, maintaining goal
... see [Gochet 2013]

In AI, “knowing how” to achieve a goal is often treated as being able to (or can) reach a goal (Situation Calculus, ATL, STIT). See two excellent surveys: [Gochet 13] and [Ågotnes, Goranko, Jamroga, Wooldridge 15]. For true know-how, Simply combining epistemic logic and the some strategy logic does not work.

Two observations inspired by the discussions in philosophy:

- Knowing how to achieve a goal may not entail that you *can* realize the goal now: a chef knows how to make cakes even when there is no sugar. The chef can make a cake, **given** all the ingredients and equipments are there.
- Even when you can win a lottery by luckily buying the right ticket, it does not mean you know how to win the lottery, since you cannot knowingly **guarantee** the result.

PLAN-BASED KNOWING HOW

The language is defined as follows:

$$\varphi ::= \top \mid p \mid \neg\varphi \mid (\varphi \wedge \varphi) \mid \text{Kh}(\varphi, \varphi)$$

$\text{Kh}(\psi, \varphi)$ reads *I know how to ensure that φ given ψ* . We define the universal modality $A\varphi$ as $\text{Kh}(\neg\varphi, \perp)$.

A model is simply a labeled transition system representing the (known) abilities of the agent: $(\mathcal{S}, \Sigma, \mathcal{R}, \mathcal{V})$ where:

- \mathcal{S} is a non-empty set of states;
- Σ is a non-empty set of actions (not in the language!);
- $\mathcal{R} : \Sigma \rightarrow 2^{\mathcal{S} \times \mathcal{S}}$ is a collection of transitions labelled by Σ ;
- $\mathcal{V} : \mathcal{S} \rightarrow 2^{\mathcal{P}}$ is a valuation function.

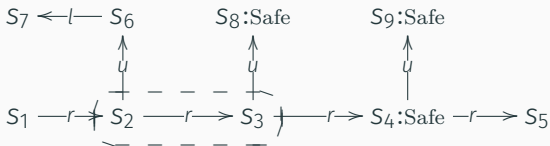
LOST WITH A MAP AT HAND



LOST WITH A MAP AT HAND



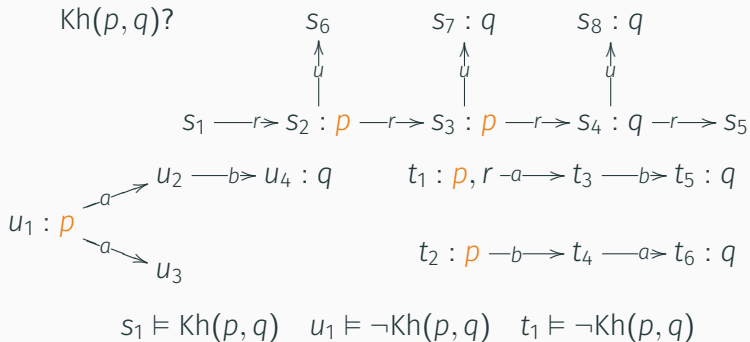
A rookie spy sneaking in an enemy building was guided by his headquarters. The communication with the HQ was lost at some point. Now someone spotted him and pulled the alarm. In panic he got lost...



Suppose he has the above map but does not know whether he is at s_2 or s_3 (the *bubble*). Does he know how to be safe? Yes! ru will make sure his safety eventually.

$\text{Kh}(\psi, \varphi)$ is true iff there is a plan σ (sequence of actions) such that you know that, given ψ , σ is always fully executable and it can get you to some φ world in the end.

$\mathcal{M}, s \models \text{Kh}(\psi, \varphi) \Leftrightarrow$ there **exists** a $\sigma \in \Sigma^*$ such that **for all** $\mathcal{M}, s' \models \psi :$
 (1) σ is **strongly executable** at s' , and
 (2) for all t if $s' \xrightarrow{\sigma} t$ then $\mathcal{M}, t \models \varphi$



$\mathcal{M}, s \models \text{A}\varphi \Leftrightarrow \text{Kh}(\neg\varphi, \perp) \Leftrightarrow$ for all $t \in \mathcal{S}, \mathcal{M}, t \models \varphi$

Achieving while maintaining [Li & Wang ICLA17]: $\text{Khm}(\psi, \chi, \varphi)$ means knowing how to achieve φ given ψ by **only passing** χ -states in-between.

$\mathcal{M}, s \models \text{Khm}(\psi, \chi, \varphi) \Leftrightarrow$ there **exists** a $\sigma \in \Sigma^*$ s.t. **for all** $\mathcal{M}, s' \models \psi$:
 (1) σ is strongly **χ** -executable at s' , and
 (2) for all t if $s' \xrightarrow{\sigma} t$ then $\mathcal{M}, t \models \varphi$

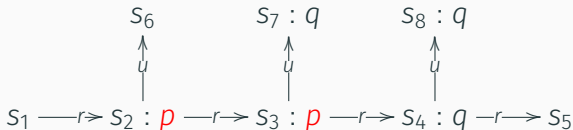
Stopping means achieving [Li Studies in Logic 17]: $\text{Khw}(\psi, \varphi)$ means knowing how to achieve φ when the execution **stops**.

$\mathcal{M}, s \models \text{Khw}(\psi, \varphi) \Leftrightarrow$ there **exists** a $\sigma \in \Sigma^*$ s.t. **for all** $\mathcal{M}, s' \models \psi$:
 for all t if $s' \xrightarrow{\sigma}_W t$ then $\mathcal{M}, t \models \varphi$

where $s \xrightarrow{\sigma}_W t$ means that the execution of σ from s may terminate at t . E.g., $\text{Khw}(p, q)$ is true below.

$$t : q \leftarrow a \mid s : p \xrightarrow{a} w : \neg a \rightarrow u : q$$

DIFFERENCES OF THE THREE KNOW-HOW OPERATORS



- $\text{Kh}(p, q)$ holds: ru is the only strongly executable witness.
- $\text{Khm}(p, p, q)$ fails: ru not only passes p states.
- $\text{Khw}(p, \neg p \wedge \neg q)$ holds, as witnessed by rrr . However, $\text{Kh}(p, \neg p \wedge \neg q)$ fails since rrr is not strongly executable.

Clearly $\text{Kh}(p, q)$ can be defined by $\text{Khm}(p, \top, q)$.

PROOF SYSTEM FOR THE FIRST SEMANTICS [WANG LOR15]

Axioms

TAUT all axioms of propositional logic

DISTU $A p \wedge A(p \rightarrow q) \rightarrow A q$

COMPKh $Kh(p, r) \wedge Kh(r, q) \rightarrow Kh(p, q)$

EMP $A(p \rightarrow q) \rightarrow Kh(p, q)$

TU $A p \rightarrow p$

4KU $Kh(p, q) \rightarrow AKh(p, q)$

5KU $\neg Kh(p, q) \rightarrow A\neg Kh(p, q)$

Rules

MP $\frac{\varphi, \varphi \rightarrow \psi}{\psi}$

NECU $\frac{\psi}{A\psi}$

SUB $\frac{\varphi(p)}{\varphi[\psi/p]}$

Provable:

PREKh: $Kh(Kh(p, q) \wedge p, q)$, **POSTKh**: $Kh(r, Kh(p, q) \wedge p) \rightarrow Kh(r, q)$

MONO: from $\vdash \varphi \rightarrow \psi$ infer $\vdash Kh(\chi, \varphi) \rightarrow Kh(\chi, \psi)$.

Axioms

TAUT all axioms of propositional logic

DISTU $Ap \wedge A(p \rightarrow q) \rightarrow Aq$

COMPKh $Kh(p, o, r) \wedge Kh(r, o, q) \wedge A(r \rightarrow o) \rightarrow Kh(p, o, q)$

EMP $A(p \rightarrow q) \rightarrow Kh(p, \perp, q)$

TU $Ap \rightarrow p$

4KU $Kh(p, o, q) \rightarrow AKh(p, o, q)$

5KU $\neg Kh(p, o, q) \rightarrow A\neg Kh(p, o, q)$

UKhm $A(p' \rightarrow p) \wedge A(o \rightarrow o') \wedge A(q \rightarrow q') \wedge Kh(p, o, q) \rightarrow Kh(p', o', q')$

OneKhm $Kh(p, o, q) \wedge \neg Kh(p, \perp, q) \rightarrow Kh(p, \perp, o)$

Rules are MP, NECU, SUB as before.

PROOF SYSTEM FOR THE THIRD SEMANTICS [LI 17]

Axioms		Rules
TAUT	all axioms of propositional logic	MP
DISTU	$Ap \wedge A(p \rightarrow q) \rightarrow Aq$	NECU
AKh	$A(p' \rightarrow p) \wedge A(q \rightarrow q') \wedge Khw(p, q) \rightarrow Khw(p', q')$	SUB
EMP	$A(p \rightarrow q) \rightarrow Khw(p, q)$	
TU	$Ap \rightarrow p$	
4KU	$Khw(p, q) \rightarrow AKhw(p, q)$	
5KU	$\neg Khw(p, q) \rightarrow A\neg Khw(p, q)$	

$s_1 : p \xrightarrow{a} s_3 : r \xrightarrow{b} s_5 : q \quad Khw(p, r) \wedge Khw(r, q) \not\vdash Khw(p, q)$

$s_2 : p, r \xrightarrow{b} s_4 : q$

EXAMPLE: CANONICAL MODEL FOR MCS Γ (FOR THE FIRST SEMANTICS)

A single canonical model does not work!

Given a maximal consistent set Γ w.r.t. SKH , let

$\Sigma_\Gamma = \{\langle \psi, \varphi \rangle \mid \text{Kh}(\psi, \varphi) \in \Gamma\}$, the canonical model for Γ is

$\mathcal{M}_\Gamma^c = \langle \mathcal{S}_\Gamma^c, \mathcal{R}^c, \mathcal{V}^c \rangle$ where:

- $\mathcal{S}_\Gamma^c = \{\Delta \mid \Delta \text{ is a MCS w.r.t. } \text{SKH} \text{ and } \Gamma|_{\text{Kh}} = \Delta|_{\text{Kh}}\}$;
- $\Delta \xrightarrow{\langle \psi, \varphi \rangle}_c \Theta$ iff $\text{Kh}(\psi, \varphi) \in \Gamma, \psi \in \Delta$, and $\varphi \in \Theta$;
- $p \in V^c(\Delta)$ iff $p \in \Delta$.

Clearly Γ is a state in \mathcal{M}_Γ^c .

Lemma (Truth lemma)

For any $\varphi \in \Gamma : \mathcal{M}_\Gamma^c, \Delta \models \varphi \iff \varphi \in \Delta$

\implies : We do not prove the contrapositive. It requires induction over the length of the witness sequence σ for the truth of $\text{Kh}(\psi, \varphi)$, where **COMPKh** plays an important role.

Theorem

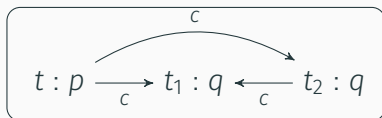
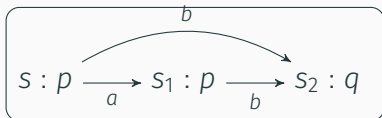
The proof systems are strongly complete w.r.t. the class of all models w.r.t the corresponding semantics.

See Yanjun Li's PhD thesis for decidability of these logics via finite canonical models.

- The $\exists\forall$ -schema in the semantics is similar to the neighborhood semantics for modal logic.
- We are inspired by *monotonic bisimulation* studied by Marc Pauly about Game Logic and Helle Hvid Hansen about monotonic neighborhood modal logic.
- In monotonic bis, wZv implies:
For any $X \in \nu_{\mathcal{M}}(w)$, there is $X' \in \nu_{\mathcal{N}}(v)$ such that for all $x' \in X'$ there is $x \in X$ such that xZx' .
- The “neighborhood” can be viewed as the collection of the sets that the agent can ensure to achieve by some plans.
- The extra complications are due to the fact that $\text{Kh}(\psi, \varphi)$ is global.

We write $U \xrightarrow{\sigma} V$ whenever σ is strongly executable for all $u \in U$, and V is the set of states reachable from U after executing σ .

We write $U \rightarrow V$ whenever there is a $\sigma \in \Sigma^*$ such that $U \xrightarrow{\sigma} V$.



- $X \xrightarrow{c} X$ thus $X \rightarrow X$ for all subsets of the state space.
- $\{s, s_1\} \xrightarrow{b} \{s_2\}$ thus $\{s, s_1\} \rightarrow \{s_2\}$.
- $\{t\} \xrightarrow{c} \{t_1, t_2\}$ thus $\{t\} \rightarrow \{t_1, t_2\}$.
- The above two models satisfy exactly the same Kh-formulas.

Definition (Kh-Bis [Fervari, Velázquez-Quesada, Wang SR17])

Let $\mathcal{M} = \langle W, \mathcal{R}, \mathcal{V} \rangle$ and $\mathcal{M}' = \langle W', \mathcal{R}', \mathcal{V}' \rangle$, a non-empty relation $Z \subseteq W \times W'$ is called an Kh-bisimulation between \mathcal{M} and \mathcal{M}' if and only if wZw' implies:

Atom: $\mathcal{V}(w) = \mathcal{V}'(w')$.

Kh-Zig: for any propositionally definable $U \subseteq W$, if $U \rightarrow V$ for some $V \subseteq W$, then there is $V' \subseteq W'$ such that

- (i) $Z[U] \rightarrow V'$ and
- (ii) for each $v' \in V'$ there is a $v \in V$ such that vZv' .

Kh-Zag: Symmetric

A-Zig: for any v in W there is a v' in W' such that vZv' .

A-Zag: Symmetric

We do need the A-Zig and A-Zag in the definition, although A is definable by Kh.

Theorem (Invariance)

Let \mathcal{M}, w and \mathcal{M}', w' be two pointed models, with $\mathcal{M} = \langle W, \mathcal{R}, \mathcal{V} \rangle$ and $\mathcal{M}' = \langle W', \mathcal{R}', \mathcal{V}' \rangle$. If $\mathcal{M}, w \Leftrightarrow_{\text{Kh}} \mathcal{M}', w'$, then $\mathcal{M}, w \equiv_{\text{Kh}} \mathcal{M}', w'$.

Theorem (Hennessy–Milner)

Let $\mathcal{M} = \langle W, \mathcal{R}, \mathcal{V} \rangle$, $\mathcal{M}' = \langle W', \mathcal{R}', \mathcal{V}' \rangle$ be two finite models, $w \in W$ and $w' \in W'$. $\mathcal{M}, w \equiv_{\text{Kh}} \mathcal{M}', w'$ iff $\mathcal{M}, w \Leftrightarrow_{\text{Kh}} \mathcal{M}', w'$.

Definition

Let $\mathcal{M} = \langle W, \mathcal{R}, \mathcal{V} \rangle$ and $\mathcal{M}' = \langle W', \mathcal{R}', \mathcal{V}' \rangle$ be two relational models. A non-empty relation $Z \subseteq (W \times W')$ is called an Khm-bisimulation between \mathcal{M} and \mathcal{M}' if and only if wZw' implies:

Atom, U-Zig and U-Zag as before.

Khm-Zig: for any propositional definable $U \subseteq W$, if $U \xrightarrow{X} V$ for some $X, V \subseteq W$, then there are $X', V' \subseteq W'$ such that

- (i) $Z[U] \xrightarrow{X'} V'$,
- (ii) for each $x' \in X'$ there is a $x \in X$ such that xZx' ,
- (iii) for each $v' \in V'$ there is a $v \in V$ such that vZv' .

Khm-Zag: Symmetric

Definition

A non-empty relation $Z \subseteq (W \times W')$ is called an Khw-bisimulation if wZw' implies:

Atom, U-Zig and U-Zag as in before.

Khm-Zig for any propositional definable $U \subseteq W$, if $U \rightarrow_W V$ for some $V \subseteq W$, then

- (i) there is $V' \subseteq W'$ such that $Z[U] \rightarrow_{W'} V'$ and,
- (ii) for each $v' \in V'$ there is a $v \in V$ such that vZv' .

Khm-Zag Symmetric.

The corresponding invariance results and Hennessy–Milner-like theorems hold.

The logic of Khm is strictly more expressive than the logic of Kh but incomparable with the logic of Khw.

Features of the logics of *knowing how* so far:

- Global knowledge
- Conditional modality
- No explicit *knowing that* operator
- Based on linear plans
- No observation during plan executions

What about a local notion with explicit know-that operator?

STRATEGY-BASED KNOWING HOW

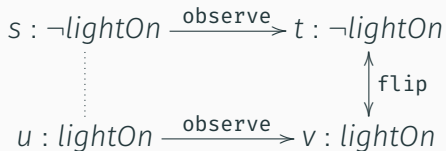
$$\varphi := p \mid \neg\varphi \mid (\varphi \wedge \varphi) \mid K\varphi \mid K_h\varphi$$

Note that we have the explicit know-that operator K in the language.

A model is a labeled transition system with an epistemic relation: $\langle \mathcal{S}, \Sigma, \mathcal{R}, \sim, \mathcal{V} \rangle$ where:

- $\langle \mathcal{S}, \Sigma, \mathcal{R}, \mathcal{V} \rangle$ is a labelled transition system as before.
- $\sim \subseteq S \times S$ is an equivalence relation (bubbles everywhere).

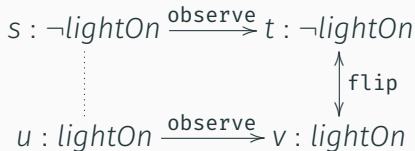
Example (reflexive arrows are omitted)



UNIFORMLY EXECUTABLE STRATEGY

- The agent's *epistemic state* at s : $[s] = \{t : s \sim t\}$
- A *strategy* is a partial function σ from epistemic states to actions.
- σ is *uniformly executable* if $\sigma([s])$ executable at every $s' \in [s]$. Empty strategy is always uniformly executable.

Example



$\sigma' = \{\{s, u\} \mapsto \text{observe}, \{t\} \mapsto \text{flip}\}$ is uniformly executable.

- $M, s \models K\varphi$ if $M, t \models \varphi$ for every t such that $s \sim t$
- $M, s \models Kh\varphi$ if there exists a *uniformly executable strategy* σ such that:
 - all complete executions starting from s terminate
 - for every final epistemic state $[t]$ after executing σ , all $t' \in [t]$ satisfies φ .

Example (M is depicted as follows)



$\sigma = \{\{s\} \mapsto a\}$ is uniformly executable, but there is an infinite execution of σ starting from s . $M, s \not\models Khp$.

A COMPLETE AXIOMATIZATION

TAUT	all axioms of propositional logic	MP	$\frac{\varphi, \varphi \rightarrow \psi}{\psi}$
DISTK	$Kp \wedge K(p \rightarrow q) \rightarrow Kq$	NECK	$\frac{\psi}{\frac{\varphi}{K\varphi}}$
T	$Kp \rightarrow p$	MonoKh	$\frac{\varphi \rightarrow \psi}{Kh\varphi \rightarrow Kh\psi}$
4	$Kp \rightarrow KKp$	SUB	$\frac{\varphi(p)}{\varphi[\psi/p]}$
5	$\neg Kp \rightarrow K\neg Kp$		
AxKtoKh	$Kp \rightarrow Khp$		
AxKhtoKhK	$Khp \rightarrow KhKp$		
AxKhtoKKh	$Khp \rightarrow KKhp$		
AxKhKh	$KhKhp \rightarrow Khp$		
AxKhbot	$\neg Kh\perp$		

- Kh is not normal

$$\not\models \text{Kh}p \wedge \text{Kh}(p \rightarrow q) \rightarrow \text{Kh}q$$

- negative introspection provable:

$$\models \neg \text{Kh}p \rightarrow \text{K}\neg \text{Kh}p$$

- sequences of modal operators reduce:

$$\models \text{K}\text{K}p \leftrightarrow \text{K}p$$

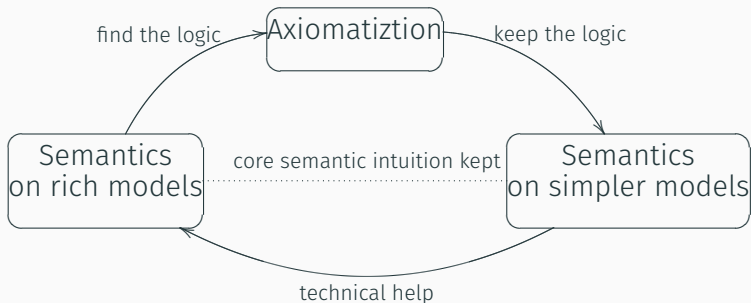
$$\models \text{Kh}\text{Kh}p \leftrightarrow \text{Kh}p$$

$$\models \text{Kh}\text{K}p \leftrightarrow \text{Kh}p$$

$$\models \text{K}\text{Kh}p \leftrightarrow \text{Kh}p$$

- sound and complete: soundness of $\text{Kh}\text{Kh}p \rightarrow \text{Kh}p$ is highly non-trivial!
- decidable: we can construct a finite canonical model.

AGAIN, WE ARE SEEKING FOR A SIMPLIFIED SEMANTICS



It also helps us to understand why it is natural to have the neighbourhood semantics in game logic and other similar logics.

The model class \mathcal{C} consists of all *mixed* epistemic models $\langle S, \sim, N, V \rangle$ satisfying the following conditions.

- For all $s \in W$, any $X, Y \subseteq W, X \in N(s)$ implies $Y \in N(s)$ (**MonoKh**).
- For all $s \in W, \emptyset \notin N(s)$ (**AxKhbot**)
- For any $s, t \in W, s \sim t$ implies $N(s) = N(t)$ (**AxKhstoKKh**).
- For all $s \in W, [s] \in N(s)$ (**AxKstoKh**).
- For all $s \in W$ and $X \subseteq W$, if $X \in N(s)$ then $Y = \{t \mid [t] \subseteq X\} \in N(s)$ (**AxKhstoKhK**)
- For all $s \in W$ and $X, Y \subseteq W$, if $X \in N(s)$, Y is definable, and $Y \in N(x)$ for all $x \in X$, we will have $Y \in N(s)$ (**AxKhKh**).

FURTHER DIRECTIONS

- Multi-agent knowing how: e.g., **one-step** coalition knowledge-how with distributed knowledge and abilities: [Naumov and Tao TARK17, AAI17]
- Goal-maintaining [Naumov and Tao AAMAS17]
- Knowingly doing [Broersen JPL2011]
- Commonly knowing how
- Comparison with various semantics of epistemic ATL
- Characterization theorems
- Logical omniscience of knowing how
- Update of *knowing how*
- Epistemic planning [Li, Yu, Wang JLC18]