

Prolog et Traitement Automatique des Langues

Éric de la Clergerie

`Eric.De_La_Clergerie@inria.fr`

ALPAGE – INRIA

`http://alpage.inria.fr`

Cours M2 LI 2008

Aujourd'hui: Représenter l'information

Cinquième partie V

Représenter l'information

La notation par termes Prolog n'est pas toujours adaptée à l'écriture **compacte**

- de termes **riches en information**
- mais en général **partiellement spécifiés** dans les productions.

Exemple d'un verbe avec ses diverses informations morpho-syntactiques :

`v(Number,Gender,Person,Tense,Mood,Aux,Voice)`

Facile de se tromper :

- sur l'arité : `v(Number,Gender,Person,Mood,Aux,Voice)`
- sur l'ordre des arguments : `v(Number,Gender,Person,Mood,Tense,Aux,Voice)`
- sur les valeurs des propriétés :
`v(pl, Gender,Person,pres,ind,Aux,Voice)`
`v(plural, Gender,Person,present,ind,Aux,Voice)`

Notation vite lourde car tout les champs doivent être explicitement remplis, même si c'est de manière anonyme : `v(,_,3,_,_,_,_)`

Mise à jour compliquée si on change l'arité d'un terme :

`v(Number,Gender,Person,Tense,Mood,Aux,Voice,SupportVerb)`

On recherche des notations :

- évitant les problèmes d'arité et d'ordre sur les propriétés
- permettant de ne fournir que les informations utiles (sous-spécification)
- permettant l'ajout de nouvelles propriétés
- évitant au maximum les erreurs en spécifiant des contraintes de bonne formation
similaire au typage en ML (CAML) ou aux DTD/schémas XML en XML
- exploitant ces contraintes pour encore plus de compacité
- exploitant ces contraintes pour plus d'efficacité pendant l'exécution

Quelques pistes

Abstraitement, un objet linguistique

- d'un certain **type**
- se caractérise par un ensemble de **propriétés** ou **traits**
- prenant leurs valeurs dans des **domaines de valeurs**, souvent finis

Les verbes de type v, avec les propriétés et valeurs

propriétés	valeurs
number	sg, pl
gender	fem, masc
person	1,2,3
tense	present, future, imperfect, past,perfect, ...
mood	indicative, subjunctive, conditional, infinitive, ...
aux	être, avoir
voice	active, passive

Proposition :

- domaines finis de valeurs
- structures de traits, typés ou non

En **DYALOG**, l'opérateur `::` permet l'unification de termes pendant la phase de lecture des programmes (plutôt que pendant la compilation ou l'exécution).

`X :: gn (Number, Gender) —> X, rel (gap (X)) .`

(presque) équivalent à :

`gn (Number, Gender) —> gn (Number, Gender), rel (gap (Number, Gender)) .`

Motivation principale pour `::` : partager des occurrences de termes complexes

Attention : Unification immédiate délicate d'utilisation :

$p(X) :- q(X :: f(a)) .$

$p(X) :- r(X :: g(b)) .$

⇒ OK

$p(X) :- q(X :: f(a)) ; r(X :: g(b)) .$

⇒ Erreur ! Échec de l'unification immédiate

Important : **Portée** de l'unification immédiate = portée des variables
= "clause" Prolog (clause, fait, production, ...)

- 1 Domaines finis
- 2 Structures de traits
- 3 Structures de traits typées

Possibilité de

- 1 nommer un ensemble fini de valeurs (atomiques)
- 2 utiliser des variables prenant leur valeur dans des sous-ensembles

```
%% Declaration
```

```
:-finite_set(mood,[cond,ger,imp,ind,inf,part,subj]).
```

```
%% Utilisation
```

```
s(Mood) —> gv(Mood::mood[inf,imp,part,ger]).
```

```
s(Mood) —> gn, gv(Mood::mood[cond,ind,subj]).
```

Spécifique **DYALOG**, mais proche :

- énumérations à la `C enum` (ou en `CAML`)
- programmation logique avec contraintes sur les domaines finis `cc(FD)`,
[Van Hentenryck](#), [V. Saraswat](#), [Y. Deville](#)
mais mise en oeuvre bien plus limitée dans **DYALOG** (pas de contraintes)

- Toutes les valeurs : `person[]` \equiv `person[1,2,3]`
- Définition par différence :

```
%% Utilisation
```

```
s(Mood) —> gv(Mood::mood[~ [cond,ind,subj]]) .  
s(Mood) —> gn, gv(Mood::mood[cond,ind,subj]) .
```

- Définition de sous-ensembles :

```
%% Déclaration
```

```
:-subset(finite, mood[cond,ind,subj]) .  
:-subset(non_finite, mood[~ [cond,ind,subj]]) .
```

```
%% Utilisation
```

```
s(Mood) —> gv(Mood::non_finite []).  
s(Mood) —> gn, gv(Mood::finite []).
```

Note : Les sous-ensembles sont du *sucre syntaxique* :

`finite [ind,subj] \rightsquigarrow mood[ind,subj]`

Un domaine fini de type T est associé à un domaine fini de valeurs \mathcal{D}_T .

- L'interprétation de FD-terme $t[a_1, \dots, a_n]$ est $(T, \{a_1, \dots, a_n\})$, avec $a_i \in \mathcal{D}_T$
(cas spécial) $\mathcal{I}(t[]) = (T, \mathcal{D}_T)$
- L'interprétation de $t[\sim[a_1, \dots, a_n]]$ est $(T, \mathcal{D}_T / \{a_1, \dots, a_n\})$
- L'unification de deux FD-termes (T, \mathcal{D}_1) et (T, \mathcal{D}_2) donne $(T, \mathcal{D}_1 \cap \mathcal{D}_2)$ avec échec de l'unification si $\mathcal{D}_1 \cap \mathcal{D}_2 = \emptyset$
- le terme (T, \mathcal{D}_1) généralise (ou **subsume**) le terme (T, \mathcal{D}_2) ssi $\mathcal{D}_2 \subset \mathcal{D}_1$

Unification : exemple

```
?-X=mood[cond, ger, imp, ind, inf],  
   Y = mood[ger, imp, ind, inf, part],  
   X=Y.
```

%% Answer:

```
X = A__12 :: mood[ger, imp, ind, inf]  
Y = A__12 :: mood[ger, imp, ind, inf]
```

Note : les deux variables X et Y ont été **redirigées** vers une nouvelle variable A_12 , liée au nouveau domaine de valeurs.

```
?- finite [] = non_finite [].  
%% => échec !
```

DYALOG propose un traitement spécial des ensembles singleton
le terme singleton $(T, \{a\})$ est considéré comme équivalent à la constante a .

?-X=mood[ind].
X = ind

Impact sur l'unification :

- l'unification de 2 FD-termes peut être une constante

?-X=mood[ind, subj], X=mood[ind, cond].
X = ind

- Un FD-terme peut être unifié avec une constante

?-X=ind, X=mood[ind, cond].
X = ind

Attention : L'interprétation des singletons est pratique mais incorrecte (formellement) !

En effet, l'unification n'est plus associative : $(t_1 \sqcup (t_2 \sqcap t_3)) \neq ((t_1 \sqcup t_2) \sqcup t_3)$

```
:-finite_set(foo,[a,b,c]).  
:-finite_set(bar,[c,d,e]).
```

```
?-X=foo[],X=bar[],X=c.
```

```
%% => Échec
```

```
?-X=foo[],X=c,X=bar[].  
    X = c
```

Les FD-termes sont utiles d'un point de vue descriptif :

$s(\text{Mood}) \longrightarrow \text{gn}, \text{gv}(\text{Mood} :: \text{finite } [])$.

équivalent à :

$s(\text{Mood}) \longrightarrow \text{gn}, \text{gv}(\text{Mood}), \{ \text{domain}(\text{Mood}, [\text{cond}, \text{ind}, \text{subj}]) \}$.

Mais FD plus efficace en pratique car évitant une exploration en retour arrière des différentes valeurs.

L'efficacité provient d'une implantation par vecteurs de bits :

```
%%          bit pos:0 1 2 3 4
:-finite_set(foo,[a,b,c,d,e]).
?-X=foo[a,c,d],
  X =.. T.    %% destructure X en T
%% Answer
X = B__5::foo[a, c, d]
T = [$SET$,B__5,foo(a, b, c, d, e),13] %% 13 = 1+4+8 = 10110
```

le DVar:::foo[a, c, d] représenté en interne par '\$SET\$(DVar,foo(a,b,c,d,e),13)

Unification \equiv **et**-booléen $\&\& : (T, bv_1) \sqcup (T, bv_2) = (T, bv_1 \&\& bv_2)$

- pas de retours arrière
- unification en temps (quasi) constant

Idem pour le test de subsumption $t_1 \sqsubseteq t_2$ ssi $bv_1 \&\& bv_2 == bv_2$

Les FD-termes sont un cas particulier de **termes déréréférençables** (**DYALOG**) :

- en Prolog, une variable X peut être liée à un terme t
 t est la valeur de X par déréréférencement
- en **DYALOG**, un deref term t_1 peut être lié à un terme t_2
 t_2 est la valeur de t_1 par déréréférencement

Les deref-termes possède une variable en 1er argument

'\$SET\$(DVar,foo(a,b,c,d,e),13)

Un deref-term est lié quand cette variable est liée.

Un dérérérencement peut suivre une chaîne de liaison

$X::foo[a,b,c,d,e] \rightsquigarrow (\text{liaison } X) Y::foo[a,b,c] \rightsquigarrow (\text{liaison } Y) b$

Très nombreux usages des domaines finis en TAL,

- morpho-syntaxe
valeurs de propriétés, comme number, gender, mood, tense,
- syntaxe (valeurs de propriétés)
- morphologie
valeurs de propriétés et classes de lettres
- phonétique, phonologie
valeurs de propriétés et classes de sons (fricative, labiale, ...)

Exemple : regexp 2 bandes

Reconnaissance d'un lemme à partir d'une forme et de règles morphologiques exprimées avec des transducteurs

```
:-finite_set (letters , [a,b,d,... , z]) .  
:-subset (voyelles , letters [a,e,i,o,u,y]) .  
:-subset (consonnes , letters [b,c,d,f,g,h,... , z]) .  
:-subset (doubles , consonnes [l,n,m]) .
```

```
suffix (v1 ,  
  ( c(X::doubles []) : c(X) ,  
    c(X) : true ,                %% même lettre doublée  
    c(Y::voyelles[e]) : c(Y) ,  %% c(e) : c(Y) aussi OK  
    true : c(consonnes[r])  
  ) ,  
  [ person: person[1,3] ,  
    tense: tense[present] ,  
    mood: mood[ind , subj]  
  ] ) .
```

```
%% appe+lle => appe+ler   1 | 3.pres.ind | subj
```

Exemple (variante)

Exploitation plus agressive de l'unification immédiate !
partage des actions plutôt que partage des lettres

```
suffix(v1,  
  ( X::c(doubles[]) : X,  
    X : true ,  
    Y::c(voyelles[e]) : Y,  
    true : c(consonnes[r])  
  ),  
  [ person: person[1,3],  
    tense: tense[present],  
    mood: mood[ind,subj]  
  ]).
```

```
%% appe+lle => appe+ler 1|3.pres.ind|subj
```

- 1 Domaines finis
- 2 Structures de traits
- 3 Structures de traits typées

Besoin naturel de représenter des ensembles de propriétés.

Plus élégant de nommer explicitement les propriétés, ce qui rend superflu :

- leur ordre d'apparition
- leur présence, quand non informative (information partielle)

⇒ même soucis dans la plupart des langages de programmation :

- `struct` en C
- `hash` en Perl (Python, Ruby, ...)
- listes de propriétés en Lisp

Notion de **structures de traits** en TAL

Notation **matricielle** des structures de traits, sous forme de **Matrice Attributs Valeurs** [AVM- *Attribute Value Matrix*]

$$\begin{bmatrix} \text{number} & \text{sg} \\ \text{person} & 2 \\ \text{tense} & \text{present} \\ \text{mood} & \text{indicative} \end{bmatrix} \equiv \begin{bmatrix} \text{person} & 2 \\ \text{mood} & \text{indicative} \\ \text{tense} & \text{present} \\ \text{number} & \text{sg} \end{bmatrix}$$

La valeur d'une propriété peut être, récursivement, une structure de traits (AVM)

$$\left[\begin{array}{l} \text{spec} \\ \text{head} \end{array} \left[\begin{array}{ll} \text{accord} & \left[\text{number} \quad \text{sg} \right] \\ \text{cat} & \text{det} \end{array} \right] \right]$$
$$\left[\begin{array}{ll} \text{accord} & \left[\text{gender} \quad \text{masc} \right] \\ \text{cat} & \text{nc} \end{array} \right]$$

Plusieurs propriétés peuvent partager une même valeur :

spec	accord	1	[number sg]
	cat	det	
head	accord	1	[gender masc]
	cat	nc	

Plusieurs propriétés peuvent partager une même valeur :

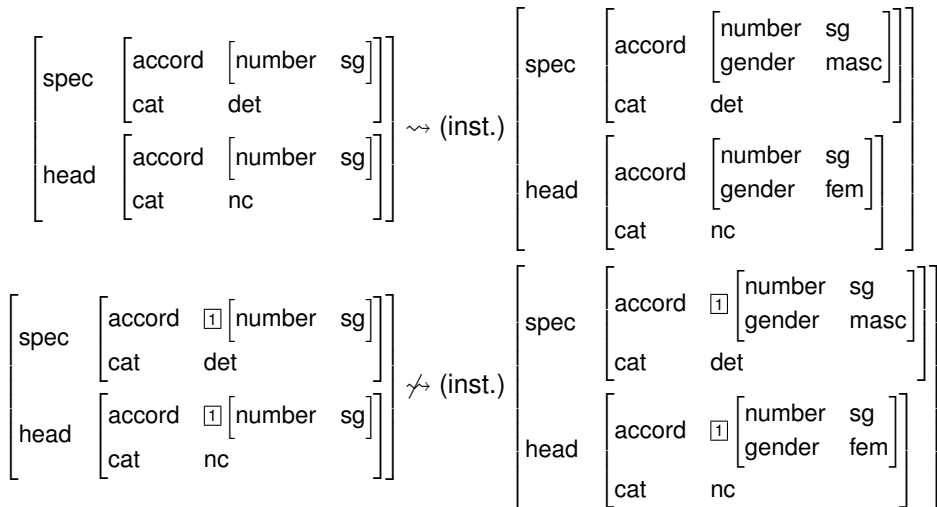
$$\begin{array}{l} \left[\begin{array}{l} \text{spec} \left[\begin{array}{l} \text{accord} \quad \boxed{1} \left[\begin{array}{l} \text{number} \quad \text{sg} \end{array} \right] \\ \text{cat} \quad \text{det} \end{array} \right] \\ \text{head} \left[\begin{array}{l} \text{accord} \quad \boxed{1} \left[\begin{array}{l} \text{gender} \quad \text{masc} \end{array} \right] \\ \text{cat} \quad \text{nc} \end{array} \right] \end{array} \right] \\ \equiv \\ \left[\begin{array}{l} \text{spec} \left[\begin{array}{l} \text{accord} \quad \boxed{1} \left[\begin{array}{l} \text{number} \quad \text{sg} \\ \text{gender} \quad \text{masc} \end{array} \right] \\ \text{cat} \quad \text{det} \end{array} \right] \\ \text{head} \left[\begin{array}{l} \text{accord} \quad \boxed{1} \left[\begin{array}{l} \text{number} \quad \text{sg} \\ \text{gender} \quad \text{masc} \end{array} \right] \\ \text{cat} \quad \text{nc} \end{array} \right] \end{array} \right] \end{array}$$

Quelques subtilités derrière la notion de réentrance :

$$\begin{array}{l} \left[\begin{array}{l} \text{spec} \left[\begin{array}{l} \text{accord} \quad \boxed{1} \left[\begin{array}{l} \text{number} \quad \text{sg} \end{array} \right] \\ \text{cat} \quad \text{det} \end{array} \right] \\ \text{head} \left[\begin{array}{l} \text{accord} \quad \boxed{1} \left[\begin{array}{l} \text{number} \quad \text{sg} \end{array} \right] \\ \text{cat} \quad \text{nc} \end{array} \right] \end{array} \right] \\ \neq \left[\begin{array}{l} \text{spec} \left[\begin{array}{l} \text{accord} \quad \left[\begin{array}{l} \text{number} \quad \text{sg} \end{array} \right] \\ \text{cat} \quad \text{det} \end{array} \right] \\ \text{head} \left[\begin{array}{l} \text{accord} \quad \left[\begin{array}{l} \text{number} \quad \text{sg} \end{array} \right] \\ \text{cat} \quad \text{nc} \end{array} \right] \end{array} \right] \end{array}$$

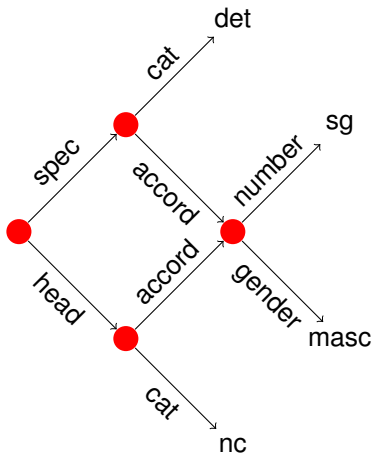
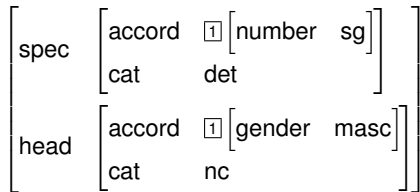
Réénontrance & instanciation

La réénontrance doit être préservée par instanciation :



Notation graphe

Jusqu'à un certain point, les structures de traits sont représentables sous formes de graphes



Note : intuitif de parler de **chemin** pour une séquence de traits dans le graphe et dans l'AVM

- Une structure de traits précise de l'information sur un objet, au travers des valeurs de traits
- L'unification de deux termes t_1 et t_2 accumule l'information $\rightsquigarrow t_1 \sqcup t_2$ (information plus spécifiée) – opération commutative et associative
- \Rightarrow algorithme :
 - ▶ tout chemin p présent dans t_1 ou t_2 est présent dans $t_1 \sqcup t_2$
 - ▶ toute valeur pour un chemin p dans t_1 ou t_2 est présente dans $t_1 \sqcup t_2$
 - ▶ tout point de partage (réentrance) dans t_1 ou t_2 est présent dans $t_1 \sqcup t_2$

- Une structure de traits précise de l'information sur un objet, au travers des valeurs de traits
- L'unification de deux termes t_1 et t_2 accumule l'information $\rightsquigarrow t_1 \sqcup t_2$ (information plus spécifiée) – opération commutative et associative
- \Rightarrow algorithme :
 - ▶ tout chemin p présent dans t_1 ou t_2 est présent dans $t_1 \sqcup t_2$
 - ▶ toute valeur pour un chemin p dans t_1 ou t_2 est présente dans $t_1 \sqcup t_2$
 - ▶ tout point de partage (réentrance) dans t_1 ou t_2 est présent dans $t_1 \sqcup t_2$

$$\left[\begin{array}{l} \text{Det} \left[\text{Accord} \left[\text{Num sing} \right] \right] \\ \text{Nom} \left[\begin{array}{l} \text{Accord} \left[\text{Num sing} \right] \\ \text{Cat N} \end{array} \right] \end{array} \right] \sqcup \left[\text{Nom} \left[\text{Accord} \left[\text{Genre Masc} \right] \right] \right]$$

- Une structure de traits précise de l'information sur un objet, au travers des valeurs de traits
- L'unification de deux termes t_1 et t_2 accumule l'information $\rightsquigarrow t_1 \sqcup t_2$ (information plus spécifiée) – opération commutative et associative
- \Rightarrow algorithme :
 - ▶ tout **chemin** p présent dans t_1 ou t_2 est présent dans $t_1 \sqcup t_2$
 - ▶ toute valeur pour un chemin p dans t_1 ou t_2 est présente dans $t_1 \sqcup t_2$
 - ▶ tout point de partage (réentrance) dans t_1 ou t_2 est présent dans $t_1 \sqcup t_2$

$$\left[\begin{array}{l} \text{Det} \left[\text{Accord} \left[\text{Num sing} \right] \right] \\ \text{Nom} \left[\begin{array}{l} \text{Accord} \left[\text{Num sing} \right] \\ \text{Cat N} \end{array} \right] \end{array} \right] \sqcup \left[\text{Nom} \left[\text{Accord} \left[\text{Genre Masc} \right] \right] \right] = \left[\begin{array}{l} \text{Det} \left[\text{Accord} \left[\begin{array}{l} \text{Num sing} \\ \text{Genre masc} \end{array} \right] \right] \\ \text{Nom} \left[\begin{array}{l} \text{Accord} \left[\text{Num sing} \right] \\ \text{Cat N} \end{array} \right] \end{array} \right]$$

Cas avec réentrance :

$$\left[\begin{array}{l} \text{Det} \left[\begin{array}{l} \text{Accord } \boxed{1} \\ \text{Num sing} \end{array} \right] \\ \text{Nom} \left[\begin{array}{l} \text{Accord } \boxed{1} \\ \text{Cat N} \end{array} \right] \end{array} \right] \sqcup \left[\text{Det} \left[\begin{array}{l} \text{Accord} \\ \text{Genre Masc} \end{array} \right] \right] = \left[\begin{array}{l} \text{Det} \left[\begin{array}{l} \text{Accord } \boxed{1} \\ \text{Num sing} \\ \text{Genre masc} \end{array} \right] \\ \text{Nom} \left[\begin{array}{l} \text{Accord } \boxed{1} \\ \text{Cat N} \end{array} \right] \end{array} \right]$$

Subsumption

- une FS t_1 généralise (ou subsume) une FS t_2 , noté $t_1 \sqsubseteq t_2$, ssi t_1 précise moins d'information que t_2
- (rappel) \sqsubseteq est un pré-ordre partiel :
 $t \sqsubseteq t$; $t_1 \sqsubseteq t_2 \wedge t_2 \sqsubseteq t_3 \Rightarrow t_1 \sqsubseteq t_3$
- \Rightarrow algorithme :
 - ▶ tout chemin dans t_1 est présent dans t_2
 - ▶ toute valeur pour un chemin p dans t_1 est présente pour p dans t_2
 - ▶ tout point de partage associé à un chemin p dans t_1 est présent dans t_2

$$\left[\begin{array}{l} \text{spec} \\ \text{head} \end{array} \left[\begin{array}{l} \text{accord} \\ \text{cat} \end{array} \left[\begin{array}{l} \text{number} \\ \text{det} \end{array} \text{sg} \right] \right] \right] \sqsubseteq \left[\begin{array}{l} \text{spec} \\ \text{head} \end{array} \left[\begin{array}{l} \text{accord} \\ \text{cat} \end{array} \left[\begin{array}{l} \boxed{1} \\ \text{det} \end{array} \left[\begin{array}{l} \text{number} \\ \text{gender} \\ \text{masc} \end{array} \right] \text{sg} \right] \right] \right] \right]$$

Les structures sont directement représentables comme des listes Prolog

$$\left[\begin{array}{l} \text{spec} \\ \text{head} \end{array} \left[\begin{array}{ll} \text{accord} & \boxed{1} \left[\text{number} \quad \text{sg} \right] \\ \text{cat} & \text{det} \end{array} \right] \right]$$
$$\left[\begin{array}{ll} \text{accord} & \boxed{1} \left[\text{gender} \quad \text{masc} \right] \\ \text{cat} & \text{nc} \end{array} \right]$$

```
?- FS = [spec: [accord: X::[number: sg],  
              cat: det],  
        head: [accord: X,  
              cat: nc]  
        ].
```

Divers problèmes d'efficacité avec la représentation ouverte en liste des FS :

- l'accès à la valeur d'un trait oblige à parcourir la list

```
fs_value(FS, F, Val) :- domain(F:Val, FS).
```

- l'unification de base Prolog n'est pas suffisante : il faut implanter une unification spécifique non efficace :

```
fs_unif(t1, t2) retourne FS
t3 := []
pour chaque f:v1 dans t1,
    si f:v2 dans t2,
        alors ajouter f:fs_unif(v1, v2) à t3
    sinon ajouter f:v1 à t3
pour chaque f:v2 dans t2,
    si f:v1 n'est_pas_dans_t1,
        alors ajouter f:v2 à t3
retourner t3
```

- \Rightarrow nombreux parcours de listes : complexité en $O(n^2)$, n taille des structures
- Plus efficace si utilisation de listes triées de traits (ou arbres balancés)
- Sinon, codage en interne (en C) appelé à partir de Prolog

Les FS ouvertes n'offrent pas de contrôle sur les traits introduits
⇒ facile de faire une erreur :

- sur le nom d'un trait : num vs number
- sur l'utilisation d'un trait dans un contexte inapproprié
ex : par unification, ajout d'un trait tense dans la FS associée à un nom

Principe :

- 1 typer les structures de traits
- 2 associer à chaque type la liste (finie) des traits **admissibles**

Cela permet :

- 1 de contrôler les occurrences des traits (noms et contextes)
- 2 d'obtenir des représentations efficaces

Déclaration des types et traits admissibles :

```
-features ([nc, gn], [number, gender]) .  
-features (v, [number, gender, person, tense, mood, voice, aux]) .  
-features (gv, [number, gender, person, mood]) .
```

Représentation :

```
s → gn{ number => Num },  
    gv{ number => Num, person => 3, mood => finite []} .  
gn{} → GN::gn{} , rel(GN) .
```


La déclaration des traits permet :

- d'associer un **rang** à chaque trait pour un type donnée
- de représenter une FS fermée par un terme Prolog standard
- dans **DYALOG**, la transformation est effectuée au moment de la lecture des FS, avec contrôle des traits :

```
DyALog> :-features(nc,[number,gender]).
```

```
DyALog> ?-X=nc{ num => sg }.
```

```
Line 2 of user_input:
```

```
    Syntax Error : feature num not allowed within feature  
    structure nc{}
```

```
DyALog> ?-X=nc{ number => sg }, X =.. T.
```

```
    X = nc{number=> sg, gender=> B__5}
```

```
    T = [nc!$ft,sg,B__5]
```

La FS $X=nc\{ \text{number} \Rightarrow \text{sg} \}$ représentée par $nc!ft(\text{sg},_)$

- les valeurs non spécifiées sont remplis par des variables
- $nc!ft$ est spécifique **DYALOG** : foncteur ft dans le domaine de noms nc

Directement fournie par l'unification standard Prolog !

```
DyALog> ?-X=nc{ number => sg }, X = nc{ gender => masc }.  
%% equiv: nc!ft(sg,_) = nc!ft(_,masc) = nc!ft(sg,masc)  
X = nc{number=> sg, gender=> masc}
```

Les structures de traits fermées sont

- moins flexibles que les FS ouvertes
- + plus sûres d'un point de vue contrôle
- + efficaces en utilisation,
- + pas de surcouche dans Prolog, (sauf dans le lecteur de termes)

Néanmoins, le contrôle n'est pas parfait !

- on peut se tromper dans la valeur d'un trait

```
?-X = v{ tense => indicative , number => plural }  
%% au lieu de: v{ mood => ind , number => pl }
```

⇒ on souhaite pouvoir préciser le domaine de valeurs d'un trait, dans le contexte d'un type.

- 1 Domaines finis
- 2 Structures de traits
- 3 Structures de traits typées**

Les **structures de traits typées** [TFS – *Typed Feature Structures*] s'appuient sur une **hiérarchie de types** :

- un type τ peut avoir des sous-types τ_1, \dots, τ_n
- un type τ peut être sous-type de plusieurs types parents (**héritage multiple**)
- un type τ peut **introduire** un trait f , en spécifiant le type $\rho_{f,\tau}$ le plus général **approprié** pour f (sous τ)
- un trait f est **approprié** pour un type τ si il est introduit par τ ou par un type ancêtre de τ

Formalisation développée par **Bob Carpenter**

- Implantation dans **ALE** en surcouche Prolog, **Gerald Penn**
<http://www.cs.toronto.edu/~gpenn/ale.html>
- autre implantation dans LKB, **Ann Copestake**
- TFS surtout utilisées pour les HSPG
- Implantation dans **DYALOG**

Fragment d'une hiérarchie pour une grammaire de génération
(*Semantic-Head-Driven Generation*, Shieber et al, dans ALE)

```
bot sub [pred, list, sem, form, agr, sign].
  form sub [finite, nonfinite].
    finite sub [].
    nonfinite sub [].
  sign sub [sentence, verbal, np, adv, p]
    intro [sem:sem].
  sentence sub [].
  verbal sub [s, vp] intro [form:form].
    s sub [].
    vp sub [] intro [subcat:subcat_list].
  np sub [det, n]
    intro [agr:agr, arg:sem].
    det sub [] intro [np_sem:sem].
    n sub [].
  adv sub [] intro [varg:sem].
  p sub [].
...

```

Notation AVM pour les TFS

```
bot sub [pred, list, sem, form, agr, sign].
sign sub [sentence, verbal, np, adv, p]
      intro [sem:sem].
sentence sub [].
verbal sub [s, vp] intro [form:form].
      s sub [].
      vp sub [] intro [subcat:subcat_list].
np sub [det, n]
      intro [agr:agr, arg:sem].
      det sub [] intro [np_sem:sem].
      n sub [].
adv sub [] intro [varg:sem].
p sub []
...

```

$$\text{np} \begin{bmatrix} \text{agr} & \text{sg1} \\ \text{arg} & \text{sem} \begin{bmatrix} \dots \end{bmatrix} \\ \text{sem} & \text{sem} \begin{bmatrix} \dots \end{bmatrix} \end{bmatrix}$$

Version puriste des TFS :
les constantes sont aussi des types
maximaux (sans sous-types) et sans traits

```
bot sub [pred, list, sem, form, agr, sign].
  agr sub [sg1, sg2, sg3, pl1, pl2, pl3].
    sg1 sub []. sg2 sub []. sg3 sub [].
    pl1 sub []. pl2 sub []. pl3 sub [].
  pred sub [decl, imp, love, call_up, leave, see, ...].
    decl sub []. imp sub [].
    leave sub []. love sub []. call_up sub []. see sub [].
  ...
  ...
```

- déclarations assez lourdes, pour des ensembles larges de constantes
- quid des ensembles ouverts, comme les entiers ou les chaînes ?

Échappements

Dans **DYALOG**, version non puriste des constantes
les constantes gérables par **échappements** pour revenir vers les types *builtins*
Prolog :

```
bot sub [humain, name, age].  
    humain intro [name: name, age: age].  
    name escape symbol.  
    age escape integer.
```

Types d'échappements : char, integer, symbol, et compound (tout terme)

Ne permet pas (encore) d'associer à un type un domaine fini de valeur

```
%% Dans un monde idéal  
bot sub [v, mood].  
    v intro [mood: mood, ...]  
    mood escape finite_set(mood)
```

À chaque type τ peut être associé la TFS la plus générale possible compatible avec les contraintes de la hiérarchie.

⇒ notion de **terme squelette** $\text{skel}(\tau)$

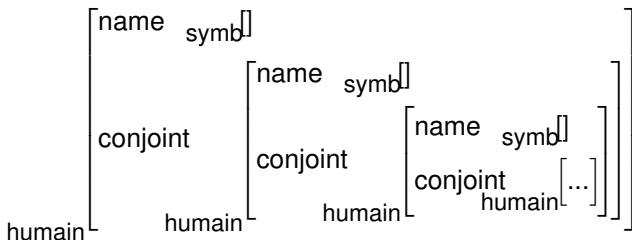
Plus précisément, si f_1, \dots, f_n sont appropriés pour τ avec les types $\rho_{f_1, \tau}, \dots, \rho_{f_n, \tau}$, alors

$$\text{skel}(\tau) = \tau \left[\begin{array}{cc} f_1 & \text{skel}(\rho_{f_1, \tau}) \\ \dots & \dots \\ f_n & \text{skel}(\rho_{f_n, \tau}) \end{array} \right]$$

Un type ne peut être directement récursif au travers d'un de ses traits,
ie $\rho_{f,\tau} \neq \tau$

```
bot sub [humain, symbol].  
  humain intro [name: symbol, conjoint: humain].  
  symbol escape symbol.
```

Motivation : conduit à un terme squelette infini



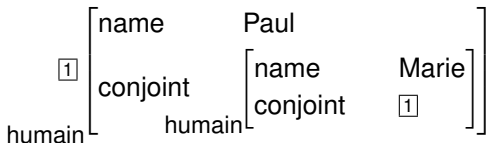
Principe : introduire des types intermédiaires sous-spécifiés

```
bot sub [humain, symbol]].
  humain sub concubin intro [name:symbol].
    concubin sub [] intro [conjoint: humain].
  symbol escape symbol.
```

Terme squelette pour concubin

```
concubin [
  [name symb]
  conjoint [name symb]
  humain ]
```

Une TFS peut être récursive et même cyclique sans que les types soient récursifs



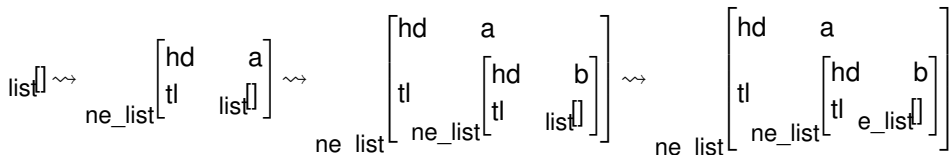
Note : la plupart des systèmes TFS refusent les termes cycliques, ils sont acceptés par **DYALOG**

On peut définir la notion de liste avec une hiérarchie de type, en utilisant un type dé-récursié `list`

```

bot sub [list, symbol, ...].
  list sub [e_list, ne_list].
    e_list sub [].
    ne_list sub []
      intro [hd:bot, tl:list].
  ...

```



Souvent, on souhaite pouvoir spécifier une liste d'éléments d'un certain type :

- non directement possible d'utiliser des listes paramétriques X list
- mais possible en sous-typant

```
list sub [e_list, ne_list, arg_list, subcat_list].
  e_list sub [].
  ne_list sub [arg_ne_list, subcat_ne_list]
    intro [hd:bot, tl:list].
  arg_list sub [e_list, arg_ne_list].
    arg_ne_list sub [] intro [hd:sem, tl:arg_list].
  subcat_list sub [e_list, subcat_ne_list].
    subcat_ne_list sub [] intro [hd:sign, tl:subcat_list].
```

Les listes sont une notion suffisamment importante pour bénéficier d'une notation spéciale dans AVM :

```
verbal sub [s, vp] intro [form:form].  
  s sub [].  
  vp sub [] intro [subcat:subcat_list].
```

$$vp \left[\begin{array}{l} \text{sem} \quad \text{sem}[] \\ \text{form} \quad \text{form}[] \\ \text{subcat} \quad \langle [\text{np}], [\text{np}], [\text{pp}] \rangle \end{array} \right]$$

- un type τ ne doit pas avoir de trait approprié f tel que $\rho_{f,\tau} = \tau$
(pas de type récursif)

- un type τ_1 peut ré-introduire un trait f avec un type ρ_{f,τ_1} (déjà introduit par un ancêtre avec le type $\rho_{f,\tau}$) mais alors ρ_{f,τ_1} doit être un sous-type de $\rho_{f,\tau}$ (instanciation des types appropriés)

```
ne_list sub [arg_ne_list, subcat_ne_list]
  intro [hd:bot, tl:list].
subcat_list sub [e_list, subcat_ne_list].
  subcat_ne_list sub [] intro [hd:sign, tl:subcat_list].
```

- si τ_1 et τ_2 introduisent le trait f ,
alors un type ancêtre τ de τ_1 et τ_2 doit introduire f
⇒ existence d'un type plus général σ_f introducteur de f Cas de hd et tl

```
list sub [e_list, ne_list, arg_list, subcat_list].
  e_list sub [].
  ne_list sub [arg_ne_list, subcat_ne_list]
    intro [hd:bot, tl:list].
  arg_list sub [e_list, arg_ne_list].
    arg_ne_list sub [] intro [hd:sem, tl:arg_list].
  subcat_list sub [e_list, subcat_ne_list].
    subcat_ne_list sub [] intro [hd:sign, tl:subcat_list].
```

L'unification entre TFS est :

- similaire à l'unification entre FS fermée
- mais doit prendre en compte l'instantiation vers des sous-types
- \Rightarrow étend strictement l'unification Prolog

Algorithme :

```
tfs_unif( $t_1$  type  $\tau_1$ ,  $t_2$  type  $\tau_2$ ) retourne TFS
 $\tau_3 :=$  mgst( $\tau_1, \tau_2$ ) ou fail
 $t_3 :=$  [ $\{\tau_3\}$ ]
pour tout  $f$  dans  $t_1$  avec  $v_1$  et  $t_2$  avec  $v_2$ ,
     $t_3 += f: v_1 \sqcap v_2 \sqcap \text{skel}(\rho_f, \tau_3)$ 
pour tout  $f$  dans  $t_1/t_2$  avec  $v_1$ 
     $t_3 += f: v_1 \sqcap \text{skel}(\rho_f, \tau_3)$ 
pour tout  $f$  dans  $t_2/t_1$  avec  $v_2$ 
     $t_3 += f: v_2 \sqcap \text{skel}(\rho_f, \tau_3)$ 
pour tout  $f$  approprié pour  $\tau_3$ , non dans  $t_1$  ou  $t_2$ 
     $t_3 += f: \text{skel}(\rho_f, \tau_3)$ 
retourne  $t_3$ 
```

L'unification entre TFS étend strictement l'unification entre termes Prolog

On peut implanter l'algorithme d'unification de manière générique en surcharge de Prolog

Mais l'information de la hiérarchie de type permet en fait de définir **hors-ligne** une opération d'unification spécialisée par paire de type τ_1 et τ_2 .

DYALOG compile les hiérarchies de types en bibliothèque C dynamique en utilisant **TFS2LIB**

Hiér. .def $\xrightarrow{\text{tfs}}$ Desc .def $\xrightarrow{\text{Perl}}$.c $\xrightarrow{\text{gcc}}$.so

Exploitation :

```
dyacc -tfs list tfs_append.pl -o tfs_append
```

TFS produit une description expansée d'une hiérarchie de types

```
TYPES = p adv n det np vp s verbal ... bot
SUBTYPES np = det n
SUBTYPES verbal = s vp
SUBTYPES sign = s vp det n sentence verbal np adv p
SUBTYPES agr = sg1 sg2 sg3 pl1 pl2 pl3
SUBTYPES form = finite nonfinite
SUBTYPES list = arg_ne_list subcat_ne_list e_list ne_list
               arg_list subcat_list
...
FEATURES = varg np_sem arg agr subcat form sem args pred tl hd
INTRO varg = adv
INTRO np_sem = det
INTRO tl = ne_list
INTRO hd = ne_list
..
```

Description expansée (skel)

```
SKEL p = sem:sem
SKEL adv = varg:sem sem:sem
SKEL n = arg:sem agr:agr sem:sem
SKEL det = np_sem:sem arg:sem agr:agr sem:sem
SKEL np = arg:sem agr:agr sem:sem
SKEL vp = subcat:subcat_list form:form sem:sem
SKEL sem = args:arg_list pred:pred
SKEL subcat_ne_list = tl:subcat_list hd:sign
SKEL subcat_list =
SKEL arg_ne_list = tl:arg_list hd:sem
SKEL arg_list =
SKEL ne_list = tl:list hd:bot
...
```



```
UNIF bot list = list
UNIF bot e_list = e_list
UNIF sign vp = vp
+      UNIF sem LEFT 1 RIGHT 3 DIRECT 3
+      INHERIT form RIGHT 2 DIRECT 2
+      INHERIT subcat RIGHT 1 DIRECT 1
...

```

Note :

- les fonctions exploitent le fait qu'un rang peut être associé à chaque trait pour un type donné.
- des fonctions de subsumption sont aussi engendrées

- Même représentation que pour les FS fermées, pour les types maximaux
- Passage en deref-termes avec un argument supplémentaire pour les termes non maximaux permet l'instantiation de types

Les informations de la hiérarchie permettent des inférences de termes à partir de chemins de traits.

```
append( e_list {}, Y, Y ).
```

```
append( A :: tl=>X, Y :: list {}, B :: tl=>Z ) :-  
  A .> hd . = B .> hd ,  
  append(X, Y, Z) .
```