# *Semantic Subtyping*

Alain Frisch (ENS Paris)

Giuseppe Castagna (ENS Paris)

Véronique Benzaken (LRI U Paris Sud)

`http://www.cduce.org/`

# *CDuce*

▷ A functional language adapted to XML applications:
  `http://www.cduce.org/`

▷ Types are pervasives in ℂDuce:

  ▷ Soundness of transformations

  ▷ Informative error messages

  ▷ Type-driven semantics, pattern matching

  ▷ Type-driven compilation and optimization

▷ This talk: theoretical foundation of ℂDuce type system.

▷ Most of the technical difficulties in the subtyping relation, which is kept simple by using a semantic approach.

# *Semantic and syntactic subtyping*

How to define a subtyping relation ?

▷ Syntactic approach: a formal system (axioms,rules)

▷ Semantic approach:

  ▷ start with a denotational model of the language

  ▷ interpret types as subsets of the model

  ▷ define subtyping as inclusion of denotations

  ▷ derive a subtyping algorithm

# *Advantages of semantic subtyping*

▷ subtyping is complete w.r.t the intuitive interpretation

▷ when $t \leq s$ does not hold, it is possible to exhibit an element of the model in the interpretation of $t$ and not of $s$ ($\rightarrow$ error message)

▷ modularization of proofs: use the semantic interpretation, not the rules of the subtyping algorithm

▷ properties "for free": transitivity of $\leq$, …

▷ Whenever possible, avoid *ad hoc* rules, formulas and algorithms: derive them from computations.

# *XDuce*

(H. Hosoya, B. Pierce, J. Vouillon)

  ▷ XDuce: a typed programming language for XML applications
  - ▷ value = XML document (tree)
  - ▷ type = regular tree language
  - ▷ subtyping = inclusion of languages

  ▷ A powerful pattern-matching operation
  - ▷ Type-driven semantics
  - ▷ Recursive patterns to extract information in the middle of a document
  - ▷ *Exact* propagation of the type from the matched expression to binding variables

# *Semantic subtyping in XDuce*

▷ Semantic subtyping: easy because no first class function, so typing of values does not depend on subtyping !

▷ Start from the typing judgment $\vdash v : t$

▷ Interpretation $[\![t]\!] = \{v \mid \ \vdash v : t\}$

▷ Subtyping $t \leq s \iff [\![t]\!] \subseteq [\![s]\!]$

# *XDuce + higher-order* $\Rightarrow \mathbb{C}$*Duce*

▷ Integrate XDuce features to a general higher order functional language

▷ Extend XDuce semantic approach ?

　▷ $[\![ t \to s ]\!] = \{ \lambda x.M \mid\, \vdash \lambda x.M : t \to s \}$

　▷ But: typing of values now depends on subtyping !

▷ Classical semantic approach ?

　▷ Define an untyped denotational model
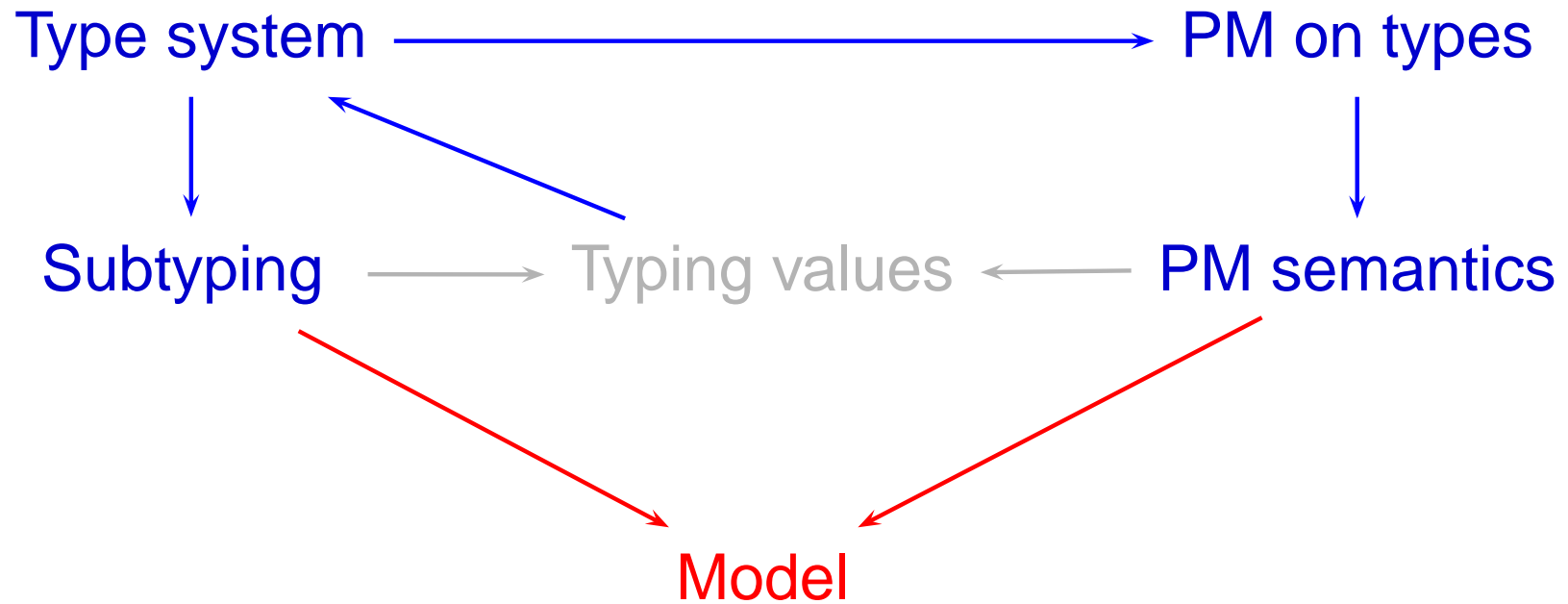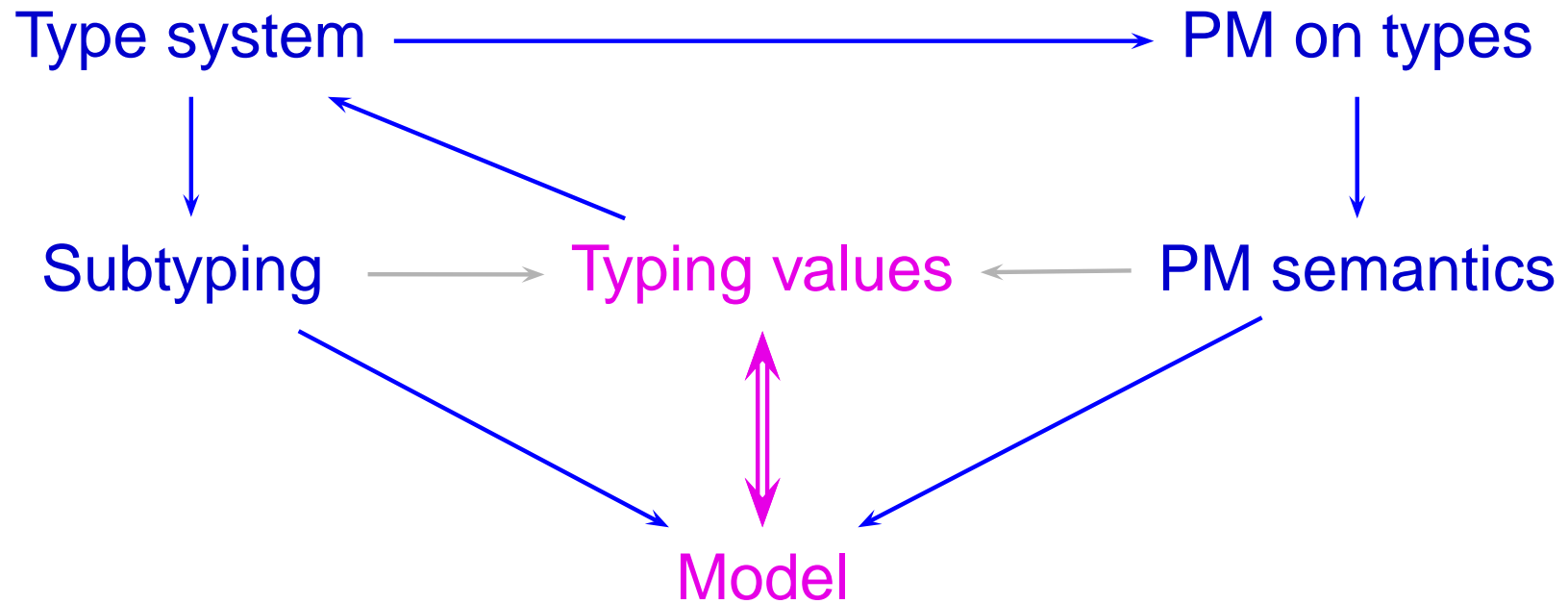
　▷ But: semantics depends on typing !

# *Circularities*

Type system $\longrightarrow$ PM on types

Subtyping $\longrightarrow$ Typing values $\longleftarrow$ PM semantics
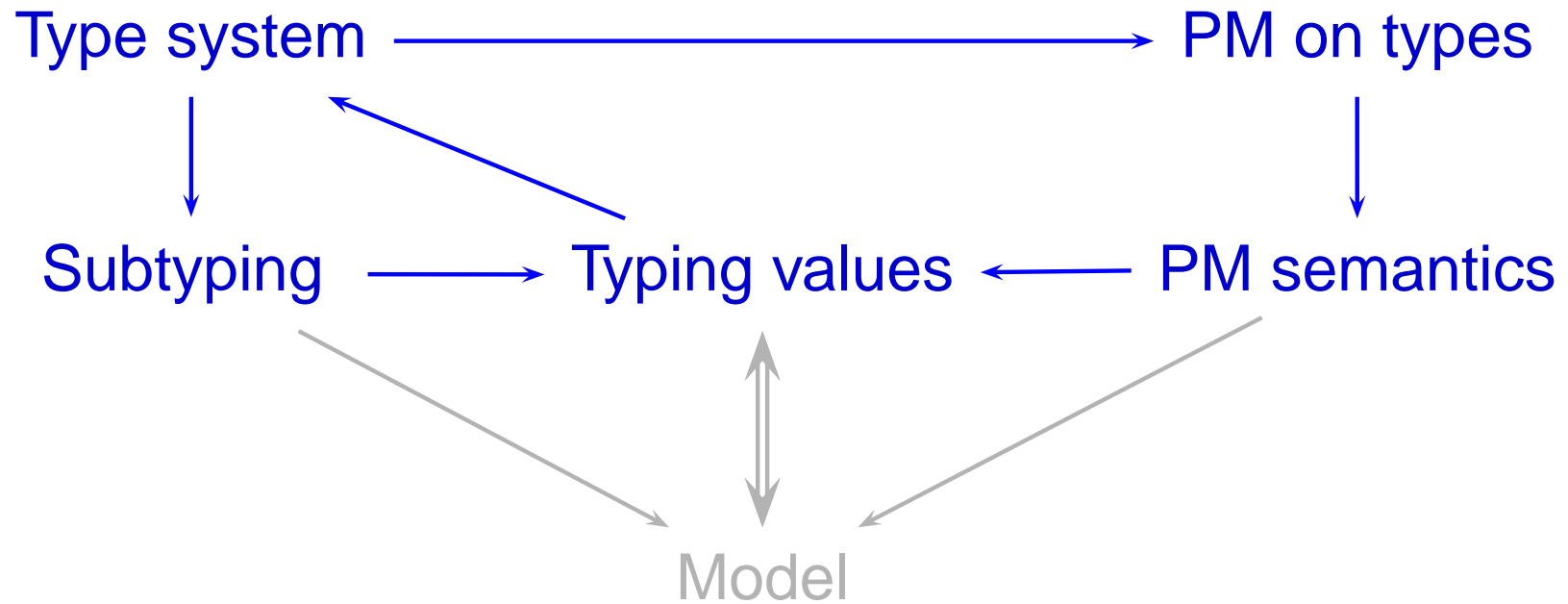
# *Circularities*

# *Circularities*

# *Circularities*

# *Types*

# *Type algebra*

$$t \quad ::= \quad b \mid t \rightarrow t \mid t \times t$$

Types:

▷ Constructors: basic, product, arrow types.

# *Type algebra*

$$
\begin{aligned}
t \quad ::= \quad & b \mid t \rightarrow t \mid t \times t \\
& \mid \quad \neg t \mid t \vee t \mid t \wedge t \mid \mathbf{0} \mid \mathbf{1}
\end{aligned}
$$

Types:

▷ Constructors: basic, product, arrow types.

▷ Arbitrary finite boolean combinations.

# *Type algebra*

$$t \quad ::= \quad b \mid t \rightarrow t \mid t \times t$$
$$\mid \quad \neg t \mid t \vee t \mid t \wedge t \mid \mathbf{0} \mid \mathbf{1}$$
$$\mid \quad \alpha \mid \mu\alpha.t$$

Types:

▷ Constructors: basic, product, arrow types.

▷ Arbitrary finite boolean combinations.

▷ Guarded recursive types.

# Subtyping and models

▷ We want to define $\leq$ by: $\boxed{t \leq s \Leftrightarrow [\![t]\!] \subseteq [\![s]\!]}$

▷ For each type, $[\![t]\!]$ is subset of a structure $\mathscr{D}$:

$$\mathscr{D} = \mathscr{D}_{\textbf{basic}} + \mathscr{D}_{\textbf{prod}} + \mathscr{D}_{\textbf{fun}} \qquad \text{with:} \qquad \mathscr{D}_{\textbf{prod}} \simeq \mathscr{D}^2$$

▷ The interpretation function $[\![\_]\!]$ must satisfy natural conditions:

$$
\begin{aligned}
[\![t_1 \times t_2]\!] &= [\![t_1]\!] \times [\![t_2]\!] \subseteq \mathscr{D}_{\textbf{prod}} \\
[\![t_1 \vee t_2]\!] &= [\![t_1]\!] \cup [\![t_2]\!] \\
[\![t_1 \wedge t_2]\!] &= [\![t_1]\!] \cap [\![t_2]\!] \\
[\![\neg t]\!] &= \mathscr{D} \setminus [\![t]\!] \\
[\![\mathbf{0}]\!] &= \varnothing
\end{aligned}
$$

# *Model condition for arrows*

$$\bigwedge_{i \in I} t_i \to s_i \leq \bigvee_{j \in J} t'_j \to s'_j$$

$$\Longleftrightarrow$$

$$(\exists j \in J)(t'_j \leq \bigvee t_i) \wedge (\forall I' \subseteq I)(t'_j \leq \bigvee_{i \in I \setminus I'} t_i) \vee (\bigwedge_{i \in I'} s_i \leq s'_j)$$

A condition on $[\![\to]\!]$ motivated by:

▷ Intuition: extensive view of functions (=binary relations)

▷ Derived from simple set-theoretic computations

▷ Makes the safety proof work

▷ The type system induces a model of values

▷ (Turns out to hold in the revelance logic B+; M. Dezani)

# *Is there at least one model !?*

▷ Yes, this one:

$$d \quad ::= \quad c \qquad\qquad\qquad\qquad c \in \mathscr{D}_{\textbf{basic}}$$
$$| \quad (d_1, d_2)$$
$$| \quad \{(d_1, d'_1), \ldots, (d_n, d'_n)\} \quad d'_i \in \mathscr{D} \cup \{\Omega\}$$

▷ It is universal (largest subtyping relation).

▷ For this model, we derive an algorithm to compute $\leq$.

# *The language*

# *The language: syntax*

$$
\begin{array}{llll}
e & ::= & x & \text{variable} \\
& | & \boldsymbol{\mu} f^{(t_1 \to s_1; \, \cdots \, ; t_n \to s_n)}(x).e & \text{abstraction} \\
& | & e_1 e_2 & \text{application} \\
& | & c & \text{constant} \\
& | & (e_1, e_2) & \text{pair} \\
& | & \texttt{match } e \texttt{ with } p_1 \Rightarrow e_1 \,|\, p_2 \Rightarrow e_2 & \text{pattern matching}
\end{array}
$$

# *Type system*

$$\frac{}{\Gamma \vdash c : t_c} \qquad \frac{}{\Gamma \vdash x : \Gamma(x)} \qquad \frac{\Gamma \vdash e_1 : t_1 \quad \Gamma \vdash e_2 : t_2}{\Gamma \vdash (e_1, e_2) : t_1 \times t_2}$$

$$\frac{\Gamma \vdash e_1 : t \rightarrow s \quad \Gamma \vdash e_2 : t}{\Gamma \vdash e_1 e_2 : s} \qquad \frac{\Gamma \vdash e : s \leq t}{\Gamma \vdash e : t}$$

$$t \equiv \bigwedge t_i \rightarrow s_i$$

$$\frac{\Gamma, (x : t_i), (f : t) \vdash e : s_i}{\Gamma \vdash \boldsymbol{\mu} f^{(t_1 \rightarrow s_1; \ldots; t_n \rightarrow s_n)}(x).e : t}$$

# *Type system*

$$\frac{}{\Gamma \vdash c : t_c} \qquad \frac{}{\Gamma \vdash x : \Gamma(x)} \qquad \frac{\Gamma \vdash e_1 : t_1 \quad \Gamma \vdash e_2 : t_2}{\Gamma \vdash (e_1, e_2) : t_1 \times t_2}$$

$$\frac{\Gamma \vdash e_1 : t \to s \quad \Gamma \vdash e_2 : t}{\Gamma \vdash e_1 e_2 : s} \qquad \frac{\Gamma \vdash e : s \leq t}{\Gamma \vdash e : t}$$

$$\frac{t \equiv \bigwedge t_i \to s_i \quad t_\neg \equiv \bigwedge \neg(t'_j \to s'_j) \text{ with } \forall j.\ t \not\leq t'_j \to s'_j \quad \Gamma, (x : t_i), (f : t) \vdash e : s_i}{\Gamma \vdash \boldsymbol{\mu} f^{(t_1 \to s_1; \ldots; t_n \to s_n)}(x).e : t \wedge t_\neg}$$

▷ Need it for subject reduction ...

▷ ... but can discard it for typechecking.

# *Results*

▷ Classical syntactical results, such as admissibility of conjunction rule and subsumption elimination.

▷ Subject reduction for a small step semantics.

▷ A new interpretation of types as sets of values:

$$\llbracket t \rrbracket_{\mathscr{V}} = \{v \mid\ \vdash v : t\}$$

   ▷ This is indeed a model.

   ▷ It induces the same subtyping relation as the bootstrap model.

   ▷ This holds because of overloaded functions !

# *The circle is now complete*

$$\vdash v : t_1 \wedge t_2 \quad \Leftrightarrow \quad (\vdash v : t_1) \wedge (\vdash v : t_2)$$

$$\vdash v : t_1 \vee t_2 \quad \Leftrightarrow \quad (\vdash v : t_1) \vee (\vdash v : t_2)$$

$$\vdash v : \neg t \quad \Leftrightarrow \quad \neg(\vdash v : t)$$

$$\nvdash v : \mathbf{0}$$

$$t \leq s \quad \Leftrightarrow \quad \forall v.(\vdash v : t) \Rightarrow (\vdash v : s)$$

$$\Leftrightarrow \quad \forall \Gamma.\forall e.(\Gamma \vdash e : t) \Rightarrow (\Gamma \vdash e : s)$$

$$\vdash v : t_1 \bullet t_2 \quad \Leftrightarrow \quad \exists v_1, v_2.(\vdash v_1 : t_1) \wedge (\vdash v_2 : t_2) \wedge (v_1 v_2 \xrightarrow{*} v)$$

$$(t_1 \bullet t_2 = \min\{s \mid t_1 \leq t_2 \to s\})$$

# *Subtyping algorithm*

▷ As any coinductive relation, the subtyping algorithm can be expressed in an abstract way through a notion of simulation.

▷ A simulation is a subset $R$ of all types closed under rules like:

$$(\bigwedge_{i \in I} t_i \rightarrow s_i \setminus \bigvee_{j \in J} t'_j \rightarrow s'_j) \in R$$

$$\Rightarrow$$

$$(\exists j \in J)(t'_j \setminus \bigvee t_i) \in R \wedge (\forall I' \subseteq I)(t'_j \setminus \bigvee_{i \in I \setminus I'} t_i) \in R \vee (\bigwedge_{i \in I'} s_i \setminus s'_j) \in R$$

▷ The largest simulation is exactly the set of empty types.

▷ This abstract presentation separates implementation issues (such as caching) and the theoretical study of the algorithm.

# *Using set-theoretic semantics*

Simple set-theoretic facts yield a number of optimizations for the subtyping algorithm, avoiding exponential explosion in many cases. E.g.:

$$(t \times s) \backslash (t_1 \times s_1) \backslash \ldots \backslash (t_n \times s_n) \simeq$$

$$(t \wedge t_1 \times s \backslash s_1) \vee \ldots \vee (t \wedge t_n \times s \backslash s_n) \vee (t \backslash t_1 \backslash \ldots \backslash t_n \times s)$$

whenever

$$\forall i \neq j.\, t_i \wedge t_j \simeq 0$$

# *Patterns*

# *Patterns*

$$
\begin{array}{lll}
p & ::= & x & \text{capture} \\
& | & t & \text{type constraint} \\
& | & p_1 \wedge p_2 & \text{conjunction} \\
& | & p_1 | p_2 & \text{alternative} \\
& | & (p_1, p_2) & \text{pair} \\
& | & (x := c) & \text{constant} \\
& | & \mu\rho.p & \text{recursive pattern} \\
& | & \rho &
\end{array}
$$

Result of matching $d/p$: $\Omega$ (failure) or a substitution *Var*$(p) \to \mathscr{D}$.

# *Pattern-matching: static semantics*

▷ The pattern $p$ accepts a type characterized by:

$$\llbracket \wr p \smallint \rrbracket = \{d \mid d/p \neq \Omega\}$$

  ▷ <u>Thm</u>: it exists + algorithm

▷ If $t \leq \wr p \smallint$ and $x \in Var(p)$, type of the result for $x$ characterized by:

$$\llbracket (t/p)(x) \rrbracket = \{(d/p)(x) \mid d \in \llbracket t \rrbracket\}$$

  ▷ <u>Thm</u>: it exists + algorithm

# *Type system: pattern matching*

$$\frac{\Gamma \vdash e : s}{\Gamma \vdash \texttt{match } e \texttt{ with } p_1 {\Rightarrow} e_1 \mid p_2 {\Rightarrow} e_2 :}$$

# *Type system: pattern matching*

$$\frac{\Gamma \vdash e : s \leq \wr p_1 \int \vee \wr p_2 \int}{\Gamma \vdash \texttt{match } e \texttt{ with } p_1 \Rightarrow e_1 \mid p_2 \Rightarrow e_2 :}$$

▷ Exhaustivity checking
- $\wr p \int$ : set of values matched by $p$

# Type system: pattern matching

$$(s_1 \equiv s \wedge \langle p_1 \rangle, \; s_2 \equiv s \wedge \neg \langle p_1 \rangle)$$

$$\frac{\Gamma \vdash e : s \leq \langle p_1 \rangle \vee \langle p_2 \rangle}{\Gamma \vdash \texttt{match } e \texttt{ with } p_1 \Rightarrow e_1 \mid p_2 \Rightarrow e_2 :}$$

▷ Exhaustivity checking
- $\langle p \rangle$ : set of values matched by $p$

▷ Dispatching

# Type system: pattern matching

$$(s_1 \equiv s \wedge \wr p_1 \smallint, \; s_2 \equiv s \wedge \neg \wr p_1 \smallint)$$

$$\frac{\Gamma \vdash e : s \leq \wr p_1 \smallint \vee \wr p_2 \smallint \quad \Gamma, (s_i/p_i) \vdash e_i : t_i}{\Gamma \vdash \mathtt{match}\ e\ \mathtt{with}\ p_1 {\Rightarrow} e_1 \mid p_2 {\Rightarrow} e_2 :}$$

▷ Exhaustivity checking
- $\wr p \smallint$ : set of values matched by $p$

▷ Dispatching

▷ Matching
- $(t/p)$ : typing environment for variables bound in $p$

# *Type system: pattern matching*

$$(s_1 \equiv s \wedge \wr p_1 \backslash,\ s_2 \equiv s \wedge \neg \wr p_1 \backslash)$$

$$\frac{\Gamma \vdash e : s \leq \wr p_1 \backslash \vee \wr p_2 \backslash \quad \Gamma, (s_i/p_i) \vdash e_i : t_i}{\Gamma \vdash \texttt{match } e \texttt{ with } p_1 \Rightarrow e_1 \mid p_2 \Rightarrow e_2 : \bigvee_{\{i \mid s_i \not\simeq \mathbf{0}\}} t_i}$$

▷ Exhaustivity checking
- $\wr p \backslash$ : set of values matched by $p$

▷ Dispatching

▷ Matching
- $(t/p)$ : typing environment for variables bound in $p$

▷ Result
- Discard useless branches

# *Typing patterns: an example*

▷ Lists *à la* Lisp (head,tail) + terminator. Consider the following pattern and types:

$$
\begin{aligned}
p &= \mu\rho.(x \wedge t_0, \mathbf{1})|(\mathbf{1}, \rho) \\
t &= \mu\alpha.(s_1 \times (s_2 \times \alpha)) \vee \texttt{nil} \\
t' &= t \wedge \big\{p\big\}
\end{aligned}
$$

▷ Then:

$$
\begin{aligned}
&\text{If } s_1 \leq t_0: && (t'/p)(x) = s_1 \\
&\text{If } s_1 \wedge t_0 \simeq \mathbf{0}: && (t'/p)(x) = s_2 \wedge t_0 \\
&\text{Otherwise:} && (t'/p)(x) = (s_1 \vee s_2) \wedge t_0
\end{aligned}
$$

# *Pattern algorithms (1)*

$$\wr x \int \quad\quad \equiv \quad \mathbf{1}$$

$$\wr t \int \quad\quad \equiv \quad t$$

$$\wr (x := c) \int \quad \equiv \quad \mathbf{1}$$

$$\wr p_1 | p_2 \int \quad \equiv \quad \wr p_1 \int \vee \wr p_2 \int$$

$$\wr p_1 \wedge p_2 \int \quad \equiv \quad \wr p_1 \int \wedge \wr p_2 \int$$

$$\wr (p_1, p_2) \int \quad \equiv \quad \wr p_1 \int \times \wr p_2 \int$$

# Pattern algorithms (2)

$$(t'/x)(x) \equiv t'$$

$$(t'/p_1|p_2)(x) \equiv ((t' \wedge \langle p_1 \rangle)/p_1)(x) \vee ((t' \wedge \neg \langle p_1 \rangle)/p_2)(x)$$

$$(t'/p_1 \wedge p_2)(x) \equiv (t'/p_i)(x) \quad \text{if } x \in \textit{Var}(p_i)$$

$$(t'/(p_1,p_2))(x) \equiv \bigvee_{(t_1,t_2) \in \pi(t')} (t_1/p_1)(x) \times (t_2/p_2)(x)$$
$$\text{if } x \in \textit{Var}(p_1) \cap \textit{Var}(p_2)$$

$$(t'/(p_1,p_2))(x) \equiv (\pi_1(t')/p_1)(x) \quad \text{if } x \in \textit{Var}(p_1) \setminus \textit{Var}(p_2)$$

$$(t'/(p_1,p_2))(x) \equiv (\pi_2(t')/p_2)(x) \quad \text{if } x \in \textit{Var}(p_2) \setminus \textit{Var}(p_1)$$

$$(t'/(x := c))(x) \equiv t_c \quad \text{if } t' \not\simeq 0$$

$$(t'/(x := c))(x) \equiv 0 \quad \text{if } t' \simeq 0$$

# *Compiling pattern matching*

▷ Naive solution: backtracking.

▷ Less naive, but untractable: bottom-up tree automata.

▷ Adopted solution: hybrid top-down/bottom-up automata, which avoid backtracking, and take profit of static type information.

▷ Complex algorithm: set-theoretic semantics of types and patterns helps *a lot* here !