

INITIAL ALGEBRAIC SEMANTICS FOR NON CONTEXT-FREE

LANGUAGES

H.Kaphengst and H.Reichel
VEB Robotron, Zentrum für
Forschung und Technik
DDR-801 Dresden,Leningrader Str. 15

Introduction

J.A.Goguen has introduced the concept of 'Initial Algebraic Semantics' in his papers [1] and [2], and he has shown that this concept yields an elegant algebraic presentation and a simple mathematical foundation of Knuth's approach in 'Semantics of Context-free Languages',[3].

The basic idea is the fact that with every context-free grammar G there can be associated in a natural way a heterogeneous operator domain τ_G in such a way that the set of derivation trees of G forms an initial algebra in the class $M(\tau_G)$ of all τ_G -algebras. On that account a context-free programming language can be fully defined by giving a syntactical τ_G -algebra whose elements are strings over the set of terminals of G and by giving a semantical τ_G -algebra.

Both Goguen and Gladki & Dikovski [4] say that this logical simple model for defining syntax and semantics of a programming language is only applicable to context-free languages.

In the paper we will outline how the basic idea in the approach of Goguen and Knuth can be extended to any recursively enumerable language and that in this extension the usually in natural languages expressed context conditions of programming languages can be expressed in a formalized way and in a legible form. The algebraic basis for this extension is the existence of freely generated equational-partial heterogeneous algebras [5],[6].

Algebraic Basic Notions

In the following we will summarize the basic concepts of equational-partial heterogeneous algebras.

A heterogeneous operator domain is a quadruple $\tau = (B, \Omega, \delta, \epsilon)$ consisting of a finite set B of sorts, of a finite set Ω of operators, of the 'input mapping' $\delta: \Omega \rightarrow B^*$ (B^* denotes the set of strings over B, including the empty string) and the 'output mapping' $\epsilon: \Omega \rightarrow B$.

A heterogeneous algebra $A = ((A[b])_{b \in B}, (A[\omega])_{\omega \in \Omega})$ is an ordered pair of

a family $(A[b])_{b \in B}$ of sets, and of a family $(A[\omega])_{\omega \in \Omega}$ of mappings $A[\omega]: A[\delta(\omega)] \rightarrow A[\delta(\omega)]$ with $A[\delta(\omega)] = A[b_1] \times A[b_2] \times \dots \times A[b_n]$ for $\delta(\omega) = b_1 b_2 \dots b_n \in B^*$.

The set $T(\tau)$ of operator terms over τ is given by:

- (1) For every $w \in B^*$ and every integer $i \leq l(w)$ ($l(w)$ denotes the length of w) " i_w " is an operator term with $w = \delta(i_w)$ as input domain and $\delta(i_w) = w(i)$ (the i -th component of w) as output domain.
- (2) For operator terms t_1, \dots, t_n with $\delta(t_1) = \delta(t_2) = \dots = \delta(t_n)$ and for an operator $\omega \in \Omega$ with $\delta(\omega) = b_1 b_2 \dots b_n$ and $b_j = \delta(t_j)$ for $j=1, 2, \dots, n$ " $\omega(t_1, t_2, \dots, t_n)$ " is an operator term with $\delta(t_1) = \delta(t_2) = \dots = \delta(t_n) = \delta(\omega(t_1, \dots, t_n))$ as input domain and $\delta(\omega(t_1, \dots, t_n)) = \delta(\omega)$ as output domain.

An unordered pair $t_1 \equiv t_2$ of operator terms with $\delta(t_1) = \delta(t_2) = w$ and $\delta(t_1) = \delta(t_2)$ is an equation in w over τ and a set $g = \{t_j \equiv r_j \mid j \in J\}$ is an equation system in w over τ if $\delta(t_j) = \delta(r_j) = w$ for every $j \in J$. With \emptyset_w we denote for $w \in B^*$ the empty equation system in w over τ . $G(\tau)$ denotes the set of all equation systems over τ and $G_w(\tau) \subseteq G(\tau)$ denotes the subset of all equation systems in w over τ .

A partial τ -interpretation $P = ((P[b])_{b \in B}, (P[\omega])_{\omega \in \Omega})$ of a heterogeneous operator domain $\tau = (B, \Omega, \delta, \delta)$ is given by a family $(P[b])_{b \in B}$ of sets and a family $(P[\omega])_{\omega \in \Omega}$ of partial mappings

$$P[\omega]: P[\delta(\omega)] \xrightarrow{\text{onto}} P[\delta(\omega)]$$

with $P[w] = P[b_1] \times \dots \times P[b_n]$ for every $w = b_1 \dots b_n \in B^*$.

By induction we can define for every $t \in T(\tau)$ and every partial τ -interpretation P a partial mapping $P[t]: P[\delta(t)] \xrightarrow{\text{onto}} P[\delta(t)]$ in such a way that $P[t](p)$ is defined for $p \in P[\delta(t)]$ and for $t = \omega(t_1, \dots, t_n)$ iff $P[t_j](p)$ for $j=1, \dots, n$ and $P[\omega](P[t_1](p), \dots, P[t_n](p))$ are defined.

For every partial τ -interpretation P and every equation system g in w over τ there can be defined a set

$$P[g] = \bigcap_{r_j \equiv s_j \in g} P[r_j \equiv s_j]$$

with $P[r \equiv s] = \{p \in P[\delta(r)] \mid P[r](p) \text{ and } P[s](p) \text{ are both defined and equal}\}$.

An equational-partial operator domain

$$\Theta = ((B, \Omega, \delta, \delta), \pi: \Omega \rightarrow \text{Nat}, \beta: \Omega \rightarrow G(\tau))$$

is given by a heterogeneous operator domain τ , by a mapping $\pi: \Omega \rightarrow \text{Nat}$, where $\pi(\omega)$ is called the degree of partiality of ω , by

a mapping $\beta: \Omega \rightarrow G(\tau)$, where $\beta(\omega)$ is called the definition-condition of ω , such that:

- (a) $\beta(\omega)$ is an equation system in $\delta(\omega)$ over τ for every $\omega \in \Omega$;
 - (b) $\pi(\omega) = 1$ if and only if $\beta(\omega) = \emptyset_{\delta(\omega)}$;
 - (c) $\pi(\omega) = 1 + \max \{\pi(\omega') \mid \omega' \in \Omega(\beta(\omega))\}$
- $(\Omega(g)$ denotes the set of operators that are used in g).

A Θ -interpretation is a partial τ -interpretation such that for every $\omega \in \Omega$ and every $p \in P[\delta(\omega)]$ holds

$$P[\omega](p) \text{ is defined if and only if } p \in P[\beta(\omega)].$$

An ordered pair (g_1, g_2) of equation systems in w over τ is called an implication in w over Θ . Let P be a Θ -interpretation for $\Theta = (\tau, \pi, \beta)$ then we say that the implication (g_1, g_2) in $w \in B^*$ over Θ holds in P if $P[g_1] \subseteq P[g_2]$.

The implications are the natural extensions of equations in universal algebra for the above defined kind of partial heterogeneous algebras. By (\emptyset_w, g) we can express the usual equations in algebraic structures, because $P[\emptyset_w] = P[w]$ is valid for every $w \in B^*$ and every partial τ -interpretation P.

An ordered pair $S = (\Theta, \Gamma)$, where Θ is an equational-partial operator domain and Γ is a set of implications over Θ , will be called a signature. By $M(S)$ we denote the class of all Θ -interpretations P such that every implication $(g_1, g_2) \in \Gamma$ holds in P. This kind of heterogeneous partial algebras will be called equational-partial heterogeneous algebras or "equoids". For $P \in M(S)$ we write also S-model or S-equoid.

Let P and M be S-equoids. An S-homomorphism $f: P \rightarrow M$ is given by a family $f = (f[b])_{b \in B}$ of mappings $f[b]: P[b] \rightarrow M[b]$, $b \in B$, such that the following condition (H) is satisfied for every $\omega \in \Omega$:

(H) If $p \in P[\beta(\omega)]$ then $f_{\delta(\omega)}(p) = (f[b_1](p_1), \dots, f[b_n](p_n))$ for

$\delta(\omega) = b_1 b_2 \dots b_n \in B^*$ is an element of $M[\beta(\omega)]$ and

$$f_{\delta(\omega)}(P[\omega](p)) = M[\omega](f_{\delta(\omega)}(p)).$$

It can be shown that for every S in $M(S)$ freely generated equoids exist. The domains of the defining operations of freely generated equoids are in general not empty. Some of the best known examples of equoids are small categories.

A signature S_2 is called an extension of S_1 , $S_1 \subseteq S_2$, if every sort, every operator and every implication of S_1 is also contained in S_2 , and if the mappings δ_1, π_1, β_1 are the corresponding restrictions of

$\delta_2, \beta_2, \pi_2, \rho_2$.

If $\underline{S}_1 \subseteq \underline{S}_2$ is a signature extension then we have a forgetful functor $U_{\underline{S}_2, \underline{S}_1} : M(\underline{S}_2) \rightarrow M(\underline{S}_1)$ which associates with every \underline{S}_2 -equoid M the underlying \underline{S}_1 -equoid $M \downarrow \underline{S}_1$ and with $f: M \rightarrow P$ the underlying \underline{S}_1 -homomorphism $f \downarrow \underline{S}_1: M \downarrow \underline{S}_1 \rightarrow P \downarrow \underline{S}_1$.

It can be shown that for every signature extension $\underline{S}_1 \subseteq \underline{S}_2$ the forgetful functor has a left adjoint functor $F_{\underline{S}_1, \underline{S}_2} : M(\underline{S}_1) \rightarrow M(\underline{S}_2)$. This means that relatively free equoids exist. The free category of paths over a directed graph is a good example for such a relatively free equoid.

Let $\underline{S}_1 \subseteq \underline{S}_2$ be a signature extension and let M be an \underline{S}_1 -equoid. The \underline{S}_2 -equoid $F(M)$ that is freely generated by M with respect to the canonical \underline{S}_1 -homomorphism $i_M: M \rightarrow F(M) \downarrow \underline{S}_1$ is called a free extension of M , if $i_M: M \rightarrow F(M) \downarrow \underline{S}_1$ is an isomorphism.

By this notion we can inductively define the concepts of a signature chain \underline{S} and of the model class $M(\underline{S})$ of a signature chain \underline{S} :

- (I₁) Every signature \underline{S} forms a signature chain \underline{S} (of length one) with \underline{S} as associated signature. Every \underline{S} -model is also an \underline{S} -model.
- (I₂) Let \underline{S} be a signature chain with \underline{S} as associated signature and let $\underline{S} \subseteq \underline{R}$ be a signature extension. Then $\underline{R} = \underline{S} \subseteq \underline{R}$ is a signature chain with \underline{R} as associated signature. Every \underline{R} -model M , so that $M \downarrow \underline{S}$ is an \underline{S} -model, is an \underline{R} -model. Such an extension of \underline{S} to $\underline{R} = \underline{S} \subseteq \underline{R}$ is called a proper extension of \underline{S} by \underline{R} .
- (I₃) Let \underline{S} be a signature chain with \underline{S} as associated signature and let $\underline{S} \subseteq \underline{R}$ be a signature extension. Then we get a new signature chain $\underline{R} = \underline{S} \sqsubseteq \underline{R}$ if for every \underline{S} -model M a free extension $F(M)$ of M to an \underline{R} -equoid exists. The associated signature of $\underline{R} = \underline{S} \sqsubseteq \underline{R}$ is \underline{R} , and we say that $\underline{R} = \underline{S} \sqsubseteq \underline{R}$ arises from \underline{S} through a weak extension by \underline{R} . The \underline{R} -models are exactly those \underline{R} -equoids that are free extensions of \underline{S} -models.

It is evident that not for every signature chain \underline{S} with \underline{S} as associated signature and for every signature extension $\underline{S} \subseteq \underline{R}$ the weak extension $\underline{R} = \underline{S} \sqsubseteq \underline{R}$ exists. In general it is recursively undecidable whether or not the weak extension of a signature chain exists.

The weakly added sorts and operators in a signature chain \underline{S} with \underline{S} as associated signature have to be interpreted as minimal solutions of the weakly added implications. In the case of total operations, $\pi(\omega)=1$ for all $\omega \in \underline{R}$, we get exactly the definition of abstract data types by

Goguen and others [7]. The more general concept of abstract data types in the language of signature chains can now be used to describe the context conditions of the production rules in the definition of a programming language.

Equation Grammars

If we convert the construction of a heterogeneous operator domain τ_G from a context-free grammar G , as described in [1] and [2], and if we start from the more general concept of signature chains, we get some kind of grammars, which we will call "equation grammars".

An equation grammar $G=(T, N, W, (\alpha_i)_{i=1, \dots, k})$ consists of a finite set $T = \{x_1, \dots, x_m\}$ of so-called terminals, of a finite set $N = \{X_1, \dots, X_n\}$ of so-called non-terminals with $N \cap T = \emptyset$, of a so-called root $w \in N$, and of a finite sequence $(\alpha_i)_{i=1, \dots, k}$ of production rules or attribute definitions or definitions of attribute functions.

A production rule is a triple consisting of a non-terminal, of a finite word over $T \cup N$, and of a domain condition, which is a possibly empty conjunction of term equations, where the terms are built up from attribute functions. For production rules we use the denotation

$$X_{j_0} \leftarrow w_{i_0} X_{j_1} w_{i_1} X_{j_2} \dots X_{j_k} w_{i_k} \text{ iff } \tau_1 = \tau'_1 \& \dots \& \tau_r = \tau'_r \text{ ffi}$$

where $X_j \in N, w_i \in T^*$.

If the domain condition of a production rule is the empty conjunction then we write nothing, and this means that the applicability of the concerned production rule is not restricted.

An attribute definition is a triple consisting of a finite set of attribute sorts, of a finite set of defining operators and of a finite set of defining implications. It is possible to define in one definition of attribute functions simultaneously some attribute functions by a recursive system of implications. A definition of attribute functions contains the names of attribute functions that will be defined, contains the descriptions of the domains and ranges of the defined attribute functions, and contains a finite set of implications. The range of an attribute function is either an attribute sort, defined before, or a non-terminal.

For want of room we can't give here a complete definition of equation grammars. The informal description of equation grammars will be explained by an example which describes the set of all such finite trees, whose nodes are marked with finite subsets of $\{a, b, \dots, z\}$ in such a way that for every path starting from the root of the tree the subsets corresponding to the nodes of the path are disjoint.

The equation grammar G is given by $T = \{ a, b, \dots, z, [], \{ \}, \{ , ; \} \}$, $N = \{ \text{char}, \text{char-sequ}, \text{tree}, \text{tree-sequ} \}$, $W = \text{tree}$ and the following sequence:

- $\alpha_1 : \text{char} \leftarrow a$
- $\alpha_2 : \leftarrow b$
- $\alpha_{26} : \cdot \cdot \cdot \leftarrow z$
- $\alpha_{27} : \underline{\text{attribute}}: \text{boole} \mid t, f: () \rightarrow \text{boole} \quad \underline{\text{end}}$
- $\alpha_{28} : \underline{\text{a-functions}}: \text{equal}:(\text{char}, \text{char}) \rightarrow \text{boole}$
 - with $\text{equal}(a, a) = \text{equal}(b, b) = \dots = \text{equal}(z, z) = t$,
 - $\text{equal}(a, b) = \dots = \text{equal}(y, z) = f \quad \underline{\text{end}}$
- $\alpha_{29} : \underline{\text{attribute}}: \text{set} \mid \lambda: () \rightarrow \text{set}, k:(\text{set}, \text{char}) \rightarrow \text{set},$
 - $U:(\text{set}, \text{set}) \rightarrow \text{set}, \epsilon:(\text{char}, \text{set}) \rightarrow \text{boole},$
 - $\cap:(\text{set}, \text{set}) \rightarrow \text{set}$
 - with $k(k(\text{set}, \text{char}), \text{char}) = k(\text{set}, \text{char}) \underline{\text{in}}(\text{set}, \text{char}),$
 - $k(k(\text{set}, \text{char}_1), \text{char}_2) = k(k(\text{set}, \text{char}_2), \text{char}_1)$
 - in($\text{set}, \text{char}, \text{char}$),
 - $U(\text{set}, \lambda) = \text{set}, U(\text{set}_1, \text{set}_2) = U(\text{set}_2, \text{set}_1),$
 - $U(\text{set}_1, k(\text{set}_2, \text{char})) = k(U(\text{set}_1, \text{set}_2), \text{char}) \underline{\text{in}}(\text{set}, \text{set}, \text{char}),$
 - $\epsilon(\text{char}, \lambda) = f,$
 - if $\epsilon(\text{char}_1, \text{set}) = t \text{ then } \epsilon(\text{char}_1, k(\text{set}, \text{char}_2)) = t \text{ in } (\text{char}, \text{set}, \text{char}) \underline{\text{fi}}$,
 - if $\epsilon(\text{char}_1, \text{set}) = f \text{ then } \epsilon(\text{char}_1, k(\text{set}, \text{char}_2)) = \text{equal}(\text{char}_1, \text{char}_2)$
 - in($\text{char}, \text{set}, \text{char}$) fi,
 - $\cap(\text{set}, \lambda) = \lambda, \cap(\text{set}_1, \text{set}_2) = \cap(\text{set}_2, \text{set}_1),$
 - if $\epsilon(\text{char}, \text{set}_1) = t \text{ then } \cap(\text{set}_1, k(\text{set}_2, \text{char})) = k(\cap(\text{set}_1, \text{set}_2), \text{char})$
 - in($\text{set}, \text{set}, \text{char}$) fi,
 - if $\epsilon(\text{char}, \text{set}_1) = f \text{ then } \cap(\text{set}_1, k(\text{set}_2, \text{char})) = \cap(\text{set}_1, \text{set}_2)$
 - in($\text{set}, \text{set}, \text{char}$) fi end
- $\alpha_{30} : \text{char-sequ} \leftarrow \text{char}$
- $\alpha_{31} : \leftarrow \text{char-sequ} ; \text{char} \underline{\text{iff}}$
 - $\epsilon(\text{char}, \text{members}(\text{char-sequ})) = f \underline{\text{ffi}}$
- $\alpha_{32} : \underline{\text{a-functions}}: \text{members}:(\text{char-sequ}) \rightarrow \text{set}$
 - with $\text{members}(\alpha_{30}(\text{char})) = k(\lambda, \text{char}),$
 - $\text{members}(\alpha_{31}(\text{char-sequ}, \text{char})) = k(\text{members}(\text{char-sequ}), \text{char}) \underline{\text{end}}$
- $\alpha_{33} : \text{tree} \leftarrow [\{ \text{char-sequ} \};]$
- $\alpha_{34} : \leftarrow [\{ \text{char-sequ} \}; \text{tree-sequ}] \underline{\text{iff}}$
 - $\cap(\text{members}(\text{char-sequ}), \text{members}(\text{tree-sequ})) = \lambda \underline{\text{ffi}}$
- $\alpha_{35} : \text{tree-sequ} \leftarrow \text{tree}$
- $\alpha_{36} : \leftarrow \text{tree-sequ} ; \text{tree}$

α_{37} : a-functions: members:(tree-sequ) → set,
 members:(tree) → set
with members([char-sequ];)=members(char-sequ),
 members([{char-sequ}; tree-sequ])=
 = \cup (members(char-sequ),members(tree-sequ)),
 members(tree-sequ ; tree) =
 = \cup (members(tree-sequ),members(tree)),
 members(α_{35} (tree))=members(tree) end

An example for a tree defined by G is given by

[{a;b};[{c;d;e};];[{c;f};[{g;h};];[{g;r;t};];[{d;e};]].

The example demonstrates that in analogy to BNF-grammars the production rules contain their syntactical interpretation as string functions over the set T of terminals of G.

In analogy to the approach of Knuth we can connect the semantical interpretation with every production rule of an equation grammar, and we get a complete description of syntax and semantics.

The class of formal languages definable by equation grammars is the class of typ 0 languages, i.e. recursively enumerable sets of strings.

References

- [1] Goguen,J.A.,Initial Algebraic Semantics, IEEE Conf. Rec. SWAT, 15 (1974), 63-77.
- [2] Goguen,J.A.,Semantics of Computation, Lecture Notes in Computer Science, vol. 25, 151-163.
- [3] Knuth,D.E.,Semantics of Context-free Languages, Math.Syst.Theory, 1968, 2, No 2, 127-145.
- [4] Gladki,A.V.,Dikovski,A.J.,Theory of Formal Grammars,(Russian), Probability Theory, Mathematical Statistics, Theoretical Cybernetics, vol. 10, 107-142, Moscow, 1972.
- [5] Kaphengst,H.,Reichel,H.,Operative Theorien und Kategorien von operativen Systemen, Studien zur Algebra und ihren Anwendungen, Berlin 1972, Akademie-Verlag.
- [6] Kaphengst,H.,Reichel,H., Algebraische Algorithmentheorie, VEB Robotron, Zentrum für Forschung und Technik, WIB Nr. 1, 1971.
- [7] Goguen,J.A.,Thatcher,J.W.,Wagner,E.G.,Wright,J.B., Abstract Data Types as Initial Algebras and Correctness of Data Representations, Proceedings, Conference on Computer Graphics, Pattern Recognition and Data Structures, May, 1975.