

# Context-Free Graph Grammars and Concatenation of Graphs

Joost Engelfriet \* and Jan Joris Vereijken \*\*

Department of Computer Science, Leiden University  
P.O.Box 9512, NL-2300 RA Leiden, The Netherlands  
e-mail: [engelfri@wi.leidenuniv.nl](mailto:engelfri@wi.leidenuniv.nl)

**Abstract.** An operation of concatenation is defined for graphs. This allows strings to be viewed as expressions denoting graphs, and string languages to be interpreted as graph languages. For a class  $K$  of string languages,  $\text{Int}(K)$  is the class of all graph languages that are interpretations of languages from  $K$ . For the classes REG and LIN of regular and linear context-free languages, respectively,  $\text{Int}(\text{REG}) = \text{Int}(\text{LIN})$ .  $\text{Int}(\text{REG})$  is the smallest class of graph languages containing all singletons and closed under union, concatenation and star (of graph languages).  $\text{Int}(\text{REG})$  equals the class of graph languages generated by linear HR (= Hyperedge Replacement) grammars, and  $\text{Int}(K)$  is generated by the corresponding  $K$ -controlled grammars. Two characterizations are given of the largest class  $K'$  such that  $\text{Int}(K') = \text{Int}(K)$ . For the class CF of context-free languages,  $\text{Int}(\text{CF})$  lies properly inbetween  $\text{Int}(\text{REG})$  and the class of graph languages generated by HR grammars. The concatenation operation on graphs combines nicely with the sum operation on graphs. The class of context-free (or equational) graph languages, with respect to these two operations, is the class of graph languages generated by HR grammars.

## 1 Introduction

Context-free graph languages are generated by context-free graph grammars, which are usually graph replacement systems. One of the most popular types of context-free graph grammar is the Hyperedge Replacement System, or HR grammar (see, e.g., [Hab, HabKre, HabKV]). A completely different way of generating graphs is to select a number of graph operations, to generate a set of expressions (built from these operations), and to interpret the expressions as graphs. The set of expressions is generated by a classical context-free grammar generating strings (or more precisely, by a regular tree grammar). This way of generating graphs was introduced, for arbitrary objects rather than graphs, in [MezWri], where the generated sets of objects are called equational. For graphs in particular, this

---

\* The first author was supported by ESPRIT BRWG No.7183 COMPUGRAPH II.

\*\* The present address of the second author is Faculty of Mathematics and Computing Science, Eindhoven University of Technology, P.O.Box 513, NL-5600 MB Eindhoven, The Netherlands, e-mail: [janjoris@acm.org](mailto:janjoris@acm.org)

generation method was first investigated in [BauCou]. It is shown in [BauCou] that, for a particular collection of graph operations, this new graph generating method is equivalent with the HR grammar. Other work on the generation of graphs through graph expressions is in, e.g., [Cou2, CouER, Dre, Eng].

In this framework we investigate another, natural operation on graphs that was introduced (for “planar nets”) in [Hot1] (and which is a simple variation of the graph operations in [BauCou]). Due to its similarity to the concatenation of strings, we call it *concatenation of graphs*. Together with the sum operation of graphs (introduced for planar nets in [Hot1] and defined for graphs in [BauCou]) and all constant graphs, a collection of graph operations is obtained that is simpler than the one in [BauCou], but also has the power of the HR grammar (which is our first main result, proved in Section 4). Concatenation and sum satisfy some nice basic properties, discussed in Section 3; in particular, all graphs can be built from a small number of elementary graphs with the operations of concatenation and sum. Thus, it suffices to use these elementary graphs in the context-free grammars that generate graph expressions.

The basic laws that are satisfied by concatenation and sum of planar nets, form the basis of the theory of x-categories developed in [Hot1] (also called strict monoidal categories, see, e.g., [EhrKKK, Ben]). Free x-categories model the sets of derivation graphs of Chomsky type 0 grammars (see [Hot2, Ben]). Finite automata on such graphs are considered, e.g., in [BosDW]. The idea of using concatenation and sum in graph grammars is from [HotKM], where “logic topological nets” are generated by graph grammars (with parallel rewriting). Our first main result (mentioned above) confirms the naturalness of these operations.

Our main interest in this paper is in the generation of graphs through graph expressions that use concatenation only. Since graph concatenation is associative, an expression that is built from constant graphs by concatenation, is essentially the same as a string. This shows that we can use arbitrary context-free grammars as graph grammars, by just interpreting the generated strings as graphs. More generally, every class  $K$  of string languages determines a class  $\text{Int}(K)$  of graph languages:  $\text{Int}(K)$  is the set of all graph languages  $h(L)$  where  $h$  is an “interpretation” and  $L$  is a string language from  $K$ . An *interpretation* of an alphabet  $A$  is a mapping  $h$  that associates a graph  $h(a)$  with every symbol  $a$ ; it is extended to strings over  $A$  by  $h(a_1 \cdots a_n) = h(a_1) \circ \cdots \circ h(a_n)$ , where  $\circ$  denotes concatenation of graphs. Thus, symbols are interpreted as graphs, strings are interpreted as graphs (by interpreting string concatenation as graph concatenation), and string languages are interpreted as graph languages. Note that an interpretation looks like a semi-group homomorphism; however, it is not exactly one, because concatenation on graphs is, in fact, a partial operation. More precisely, graphs are typed, and concatenation is defined only if the types “fit”. In fact, as in [Hab], our graphs are equipped with a designated sequence of “begin nodes” and a designated sequence of “end nodes” (generalizing the idea that strings have a beginning and an end). A graph  $g_1$  can be concatenated with a graph  $g_2$  only if the sequence of end nodes of  $g_1$  has the same length as the sequence of begin nodes of  $g_2$ . Their concatenation  $g_1 \circ g_2$  is obtained by identifying

each end node of  $g_1$  with the corresponding begin node of  $g_2$  (just as strings are concatenated by identifying the end of the first string with the beginning of the second).

We investigate  $\text{Int}(K)$  for specific  $K$  (such as the class REG of regular languages, the class CF of context-free languages, and the class LIN of linear context-free languages), but also for arbitrary  $K$  (satisfying some mild closure properties). In Section 5, after defining the notion of interpretation, we show that the graph languages in  $\text{Int}(\text{REG})$  are exactly those that can be denoted by regular expressions, built from singleton graph languages with the operations of union, concatenation, and star (on graph languages). We also show that  $\text{Int}(\text{REG}) = \text{Int}(\text{LIN})$  and that it equals the class LIN-HR of graph languages generated by *linear* HR grammars. This suggests that regularity and linearity are the same for graph languages. The class  $\text{Int}(\text{CF})$  contains, as expected, exactly those graph languages that can be generated by expression generating context-free grammars that do not use the sum operation. Thus, by our first main result, it is included in the class of graph languages generated by HR grammars. The inclusion is proper, due to the close connection between graph concatenation and the pathwidth of graphs: every graph language in  $\text{Int}(K)$  is of bounded pathwidth (and graph languages of unbounded pathwidth, such as the set of trees, can be generated by HR grammars).

Generalizing the result that  $\text{Int}(\text{REG}) = \text{LIN-HR}$ , we show in Section 6 that (under the rather weak assumption that  $K$  is closed under sequential machine mappings)  $\text{Int}(K)$  is equal to  $\text{LIN-HR}(K)$ , the class of graph languages that are generated by linear HR grammars with a control language from  $K$  (with the usual notion of control).

As observed above,  $\text{Int}(\text{REG}) = \text{Int}(\text{LIN})$ . In Section 7 we investigate the question, for given  $K$  and  $K'$ , whether or not  $\text{Int}(K') = \text{Int}(K)$  (where we assume that  $K$  and  $K'$  are closed under sequential machine mappings). Trivially, for every  $K$  there is a largest class  $\bar{K}$  such that  $\text{Int}(\bar{K}) = \text{Int}(K)$ . We call this class the *extension* of  $K$ , denoted  $\text{Ext}(K)$ . Clearly, the question  $\text{Int}(K') = \text{Int}(K)$  is now reduced to the question  $\text{Ext}(K') = \text{Ext}(K)$ , which concerns classes of string languages rather than graph languages. The main result of this section is that  $\text{Ext}(K)$  consists exactly of all string languages that are in  $\text{Int}(K)$ , coding strings as graphs in the obvious way (viz., as edge-labeled chain graphs). Using the characterization in Section 6, and generalizing a result concerning the string generating power of linear HR grammars from [EngHey], we show that  $\text{Ext}(K) = 2\text{DGSM}(K)$ , the class of all languages that are images of languages from  $K$  under 2-way deterministic gsm mappings. Thus,  $\text{Int}(K') = \text{Int}(K)$  iff  $2\text{DGSM}(K') = 2\text{DGSM}(K)$ , a purely formal language-theoretic question. By the well-known result that  $2\text{DGSM}(\text{REG})$  is properly included in  $2\text{DGSM}(\text{CF})$ , we conclude that  $\text{Int}(\text{REG})$  is properly included in  $\text{Int}(\text{CF})$ .

A preliminary version of this paper was presented at the 5th International Workshop on Graph Grammars and their Application to Computer Science [EngVer]. The work is based on the Master's Thesis of the second author [Ver].

## 2 Preliminaries

### 2.1 Strings

We assume the reader to be familiar with formal language theory (see, e.g., [Ber, HopUll, Sal]). Here we just recall some of the concepts to be used.

$\mathbf{N} = \{0, 1, 2, \dots\}$  denotes the set of natural numbers. For a set  $V$ ,  $V^*$  denotes the set of all finite sequences (or strings) of elements of  $V$ . A sequence  $\langle v_1, v_2, \dots, v_n \rangle \in V^*$ , with  $v_i \in V$ , is also written as  $v_1 v_2 \cdots v_n$  (as  $\lambda$  if  $n = 0$ ). The length of a string  $w \in V^*$  is denoted  $|w|$ , and, for  $1 \leq i \leq |w|$ , its  $i$ th element is denoted  $w(i)$ . Thus, if  $w = v_1 \cdots v_n$ , then  $|w| = n$  and  $w(i) = v_i$ . Concatenation of strings is defined in the usual way.

A context-free grammar is a tuple  $G = (N, T, P, S)$  where  $N$  is the nonterminal alphabet,  $T$  is the terminal alphabet (disjoint with  $N$ ),  $P$  is the set of productions (of the form  $X \rightarrow \alpha$ , with  $\alpha \in (N \cup T)^*$ ), and  $S$  is the initial nonterminal. The language  $L(G) \subseteq T^*$  generated by  $G$  is defined in the usual way. A context-free grammar is linear if there is at most one nonterminal occurrence in each right-hand side of a production, and it is right-linear if each production is of the form  $X \rightarrow aY$  or  $X \rightarrow a$ , with  $X, Y \in N$  and  $a \in T$ . The class of languages generated by all (all linear) context-free grammars is denoted CF (LIN, respectively). By REG we denote the class of regular languages. Note that the right-linear context-free grammars generate the class of  $\lambda$ -free regular languages (i.e., those regular languages that do not contain  $\lambda$ ).

### 2.2 Graphs and graph replacement

We consider the multi-pointed, directed, edge-labeled hypergraphs of [Hab]. Such a hypergraph consists of a set of nodes and a set of (hyper)edges, just as an ordinary graph, except that an edge may have any number of sources and any number of targets, rather than just one source and one target. Each edge is labeled with a symbol from a “doubly-ranked” alphabet, in such a way that the first (second) rank of its label equals the number of its sources (targets, respectively). Finally, every hypergraph is multi-pointed in the sense that it has a designated sequence of “begin nodes”, and a designated sequence of “end nodes”; these can be used conveniently for gluing hypergraphs to each other.

Formally, a *typed* (or doubly ranked) *alphabet* is an alphabet  $\Sigma$  together with a mapping  $\text{type} : \Sigma \rightarrow \mathbf{N} \times \mathbf{N}$ . A *multi-pointed hypergraph* over  $\Sigma$  is a tuple  $g = (V, E, s, t, l, \text{begin}, \text{end})$ , where  $V$  is the finite set of nodes,  $E$  is the finite set of (hyper)edges,  $s : E \rightarrow V^*$  is the source function,  $t : E \rightarrow V^*$  is the target function,  $l : E \rightarrow \Sigma$  is the labeling function such that  $\text{type}(l(e)) = (|s(e)|, |t(e)|)$  for every  $e \in E$ ,  $\text{begin} \in V^*$  is the sequence of begin nodes, and  $\text{end} \in V^*$  is the sequence of end nodes.

For a given multi-pointed hypergraph  $g$ , its components will also be denoted by  $V_g, E_g, s_g, t_g, l_g, \text{begin}(g)$ , and  $\text{end}(g)$ . If  $|\text{begin}(g)| = m$  and  $|\text{end}(g)| = n$ , then  $g$  is said to be of *type*  $(m, n)$  and we write  $\text{type}(g) = (m, n)$ . Similarly, for an edge  $e$  of  $g$ , we write  $\text{type}(e)$  to denote  $\text{type}(l(e))$ ; thus, by the above

requirement, if  $\text{type}(e) = (m, n)$ , then  $e$  has  $m$  sources and  $n$  targets. If a multi-pointed hypergraph is of type  $(0, 0)$  and all its edges are of type  $(1, 1)$ , then it is an ordinary directed graph (with labeled edges).

For a typed symbol  $\sigma$ , with  $\text{type}(\sigma) = (m, n)$ , we denote by  $\text{atom}(\sigma)$  the multi-pointed hypergraph  $g$  of type  $(m, n)$  such that  $V_g = \{x_1, \dots, x_m, y_1, \dots, y_n\}$ ,  $E_g = \{e\}$  with  $l(e) = \sigma$ , and  $\text{begin}(g) = s(e) = \langle x_1, \dots, x_m \rangle$ , and  $\text{end}(g) = t(e) = \langle y_1, \dots, y_n \rangle$ . A multi-pointed hypergraph of the form  $\text{atom}(\sigma)$  will be called an *atom* (it is called a handle in [Hab]).

Two multi-pointed hypergraphs  $g$  and  $h$  are *disjoint* if  $V_g \cap V_h = \emptyset$  and  $E_g \cap E_h = \emptyset$ .

*From now on we will just say graph instead of multi-pointed hypergraph.* As usual we consider both concrete and abstract graphs, where an abstract graph is an equivalence class of isomorphic concrete graphs. The isomorphisms between graphs  $g$  and  $h$  are the usual ones, which, additionally, should map  $\text{begin}(g)$  to  $\text{begin}(h)$ , and  $\text{end}(g)$  to  $\text{end}(h)$ . In particular, isomorphic graphs have the same type. We are only interested in abstract graphs; concrete graphs are just used as representatives of abstract graphs. The set of abstract graphs over a typed alphabet  $\Sigma$  will be denoted  $\text{GR}(\Sigma)$ , and  $\text{GR}$  denotes the union of all  $\text{GR}(\Sigma)$  (where  $\Sigma$  is taken from some fixed, infinite set of symbols). A (typed) *graph language* is a subset  $L$  of  $\text{GR}(\Sigma)$ , for some  $\Sigma$ , such that all graphs in  $L$  have the same type  $(m, n)$ , also called the type of  $L$ , and denoted by  $\text{type}(L) = (m, n)$ .

A basic operation on graphs is the substitution of a graph for an edge (see [Hab, BauCou]). To define it formally, it is convenient to use an operation of *node identification* (or “gluing”), as follows.

Let  $g$  be a graph, and let  $R \subseteq V_g \times V_g$ . Intuitively, we wish to identify nodes  $x$  and  $y$ , for every pair  $(x, y) \in R$ . For  $x \in V_g$ , let  $[x]_R$  denote the equivalence class of  $x$  with respect to the smallest equivalence relation on  $V_g$  containing  $R$ . For  $V \subseteq V_g$ , let  $V/R = \{[x]_R \mid x \in V\}$ . For a sequence  $x = \langle x_1, \dots, x_n \rangle \in V_g^*$  with  $x_i \in V_g$ , let  $[x]_R = \langle [x_1]_R, \dots, [x_n]_R \rangle$ . Then we define the graph  $g/R$  by  $g/R = (V_g/R, E_g, s, t, l_g, [\text{begin}]_R, [\text{end}]_R)$  such that  $s(e) = [s_g(e)]_R$  and  $t(e) = [t_g(e)]_R$  for every  $e \in E_g$ .

Substitution of a graph for an edge is now defined as follows. Let  $g$  be a graph, let  $e$  be an edge of  $g$ , and let  $h$  be a graph such that  $\text{type}(h) = \text{type}(e) = (m, n)$ . We assume that  $g$  and  $h$  are disjoint (otherwise an isomorphic copy of  $h$  should be taken). Let  $g'$  be the graph that is obtained from  $g$  by removing  $e$  and adding  $h$  (disjointly), i.e.,  $g' = (V_g \cup V_h, (E_g - \{e\}) \cup E_h, s, t, l, \text{begin}(g), \text{end}(g))$ , where  $s(e) = s_g(e)$  for  $e \in E_g - \{e\}$  and  $s(e) = s_h(e)$  for  $e \in E_h$ , and similarly for  $t$  and  $l$ . Note that  $g'$  has the begin and end nodes of  $g$ . Then the *substitution* of  $h$  for  $e$  in  $g$ , denoted by  $g[e/h]$ , is the graph  $g'/R$  where  $R = \{(s_g(e)(i), \text{begin}(h)(i)) \mid 1 \leq i \leq m\} \cup \{(t_g(e)(i), \text{end}(h)(i)) \mid 1 \leq i \leq n\}$ . Thus, intuitively, after removing  $e$  and adding  $h$ , the  $i$ th source of  $e$  is identified with the  $i$ th begin node of  $h$ , and the  $i$ th target of  $e$  is identified with the  $i$ th end node of  $h$ . The notion of substitution defined here is not precisely the one in [Hab], but it is (the appropriate extension to the doubly ranked case of) the one in [BauCou]; however, they are equivalent from the point of view of graph generation by hyperedge replacement grammars.

In a substitution  $g[e/h]$ ,  $h$  can be taken as an abstract graph (in the sense that if  $h$  and  $h'$  are isomorphic, then so are  $g[e/h]$  and  $g[e/h']$ ); but  $g$  is necessarily concrete, because its concrete edge  $e$  is involved. To turn substitution into an operation on abstract graphs, we substitute graphs for all edges of  $g$  and let the graph  $h$  to be substituted for edge  $e$  be determined by the label of  $e$ . This leads us to a notion of substitution that generalizes the notion of homomorphism of strings (in formal language theory), and that we will call “replacement” (of edges by graphs).

Let  $\Sigma$  be a typed alphabet. A *replacement* is a mapping  $\phi : \Sigma \rightarrow \text{GR}$  such that  $\text{type}(\phi(\sigma)) = \text{type}(\sigma)$  for every  $\sigma \in \Sigma$ ; it is extended to a mapping from  $\text{GR}(\Sigma)$  to  $\text{GR}$  by defining, for  $g \in \text{GR}(\Sigma)$ ,  $\phi(g) = g[e_1/\phi(l(e_1))] \cdots [e_k/\phi(l(e_k))]$ , where  $E_g = \{e_1, \dots, e_k\}$ . Thus, every edge  $e$  of  $g$  with label  $l(e) = \sigma$  is replaced by the graph  $\phi(\sigma)$ . It is well known that this definition does not depend on the order  $e_1, \dots, e_k$  in which the edges are replaced (because substitution is *confluent*, cf. [Cou1]). It should also be clear that every replacement is an operation on abstract graphs: if  $g$  and  $g'$  are isomorphic, then so are  $\phi(g)$  and  $\phi(g')$ .

We denote the class of all replacements by  $\text{Repl}$ , and, for a class  $K$  of graph languages, we let  $\text{Repl}(K) = \{\phi(L) \mid \phi \in \text{Repl}, L \in K\}$ .

Another basic property of substitution is its *associativity* (see [Cou1]). In our present formulation it means that the composition of two replacements is again a replacement (as one would expect from a generalization of string homomorphism).

**Proposition 1.** *Repl is closed under composition.*

*Proof.* It can be shown, based on the associativity of substitution, that, for a replacement  $\phi$ ,  $\phi(g[e_1/h_1] \cdots [e_k/h_k]) = g[e_1/\phi(h_1)] \cdots [e_k/\phi(h_k)]$ , where  $E_g = \{e_1, \dots, e_k\}$ . Now let  $\Sigma_1$  and  $\Sigma_2$  be two typed alphabets. Let  $\phi_1 : \Sigma_1 \rightarrow \text{GR}(\Sigma_2)$  and  $\phi_2 : \Sigma_2 \rightarrow \text{GR}$  be two replacements. Define the replacement  $\phi : \Sigma_1 \rightarrow \text{GR}$  by:  $\phi(\sigma) = \phi_2(\phi_1(\sigma))$  for every  $\sigma \in \Sigma_1$ . Then, for a graph  $g$  with  $E_g = \{e_1, \dots, e_k\}$ ,  $\phi_2(\phi_1(g)) = \phi_2(g[e_1/\phi_1(l(e_1))] \cdots [e_k/\phi_1(l(e_k))]) = g[e_1/\phi_2(\phi_1(l(e_1)))] \cdots [e_k/\phi_2(\phi_1(l(e_k)))] = g[e_1/\phi(l(e_1))] \cdots [e_k/\phi(l(e_k))] = \phi(g)$ . This shows that  $\phi = \phi_2 \circ \phi_1$ .  $\square$

A useful elementary property of replacements is that, for every replacement  $\phi : \Sigma \rightarrow \text{GR}$  and every  $\sigma \in \Sigma$ ,  $\phi(\text{atom}(\sigma)) = \phi(\sigma)$ . Also, if  $\phi(\sigma) = \text{atom}(\sigma)$  for every  $\sigma \in \Sigma$ , then  $\phi(g) = g$  for every  $g \in \text{GR}(\Sigma)$ .

Besides replacement operations, there are two other, simpler operations on (abstract) graphs that will be useful. They only change the begin and end nodes of a graph. Let  $g$  be a graph. The *fold* of  $g$ , denoted  $\text{fold}(g)$ , is the same as  $g$ , except that  $\text{begin}(\text{fold}(g)) = \lambda$  and  $\text{end}(\text{fold}(g)) = \text{begin}(g) \cdot \text{end}(g)$ , where  $\cdot$  denotes concatenation of strings, as usual. The *backfold* of  $g$ , denoted  $\text{backfold}(g)$ , is the same as  $g$ , except that  $\text{begin}(\text{backfold}(g)) = \text{begin}(g) \cdot \text{end}(g)$  and  $\text{end}(\text{backfold}(g)) = \lambda$ .

### 2.3 Hyperedge replacement grammars

Hyperedge replacement grammars (or HR grammars) are context-free graph grammars that substitute graphs for edges. An *HR grammar* is a tuple  $G = (N, T, P, S)$  where  $N$  is a typed alphabet of nonterminals,  $T$  is a typed alphabet of terminals (disjoint with  $N$ ),  $P$  is a finite set of productions, and  $S \in N$  is the initial nonterminal. Every production in  $P$  is of the form  $X \rightarrow h$  with  $X \in N$ ,  $h \in \text{GR}(N \cup T)$ , and  $\text{type}(X) = \text{type}(h)$ ; moreover, we assume (without loss of generality) that *no two edges of  $h$  are labeled by the same nonterminal*.

Application of a production  $p = X \rightarrow h$  to a graph is defined as follows. Let  $g \in \text{GR}(N \cup T)$ , and let  $e \in E_g$ . Then  $p$  is applicable to  $e$  if  $l_g(e) = X$ , and the result of the application is the graph  $g[e/h]$ . We write  $g \Rightarrow_p g'$ , or just  $g \Rightarrow g'$ , if  $g'$  is the result of applying  $p$  to  $e$  of  $g$ , i.e., if  $g'$  is (isomorphic to)  $g[e/h]$ . As usual,  $\Rightarrow^*$  denotes the transitive reflexive closure of  $\Rightarrow$ . The *graph language generated by  $G$*  is  $L(G) = \{g \in \text{GR}(T) \mid \text{atom}(S) \Rightarrow^* g\}$ . Note that  $\text{type}(L(G)) = \text{type}(S)$ .

We denote by HR the class of graph languages generated by HR grammars. An HR grammar is *linear* if there is at most one nonterminal edge in each right-hand side of a production. We denote by LIN-HR the class of graph languages generated by linear HR grammars.

A fundamental property of HR grammars is formulated in the following “context-freeness lemma” (cf. Section II.2 of [Hab]). As shown in Lemma 2.14 of [Coul], it is based on the associativity of substitution. Due to the above assumption that a nonterminal occurs at most once in the right-hand side of a production, it can be stated in terms of replacements, as follows.

**Proposition 2.** *Let  $G = (N, T, P, S)$  be an HR grammar. Let  $X \rightarrow h$  be in  $P$ , and let  $g \in \text{GR}(T)$ . Let  $\text{lab}(h) = \{l_h(e) \mid e \in E_h\}$ . Then  $h \Rightarrow^* g$  if and only if there exists a replacement  $\phi : \text{lab}(h) \rightarrow \text{GR}(T)$  such that  $\phi(h) = g$  and  $\text{atom}(\sigma) \Rightarrow^* \phi(\sigma)$  for every  $\sigma \in \text{lab}(h)$ . Moreover, the length of the derivation  $h \Rightarrow^* g$  equals the sum of the lengths of all derivations  $\text{atom}(\sigma) \Rightarrow^* \phi(\sigma)$ .*

## 3 Concatenation and Sum

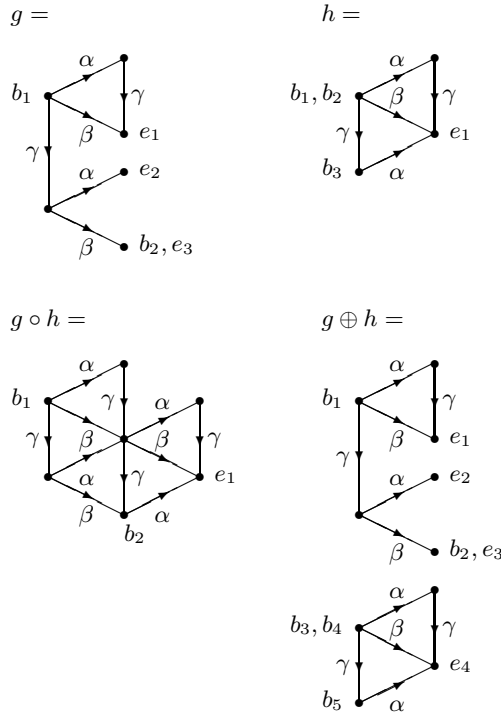
In this section we define the graph operation of concatenation, and investigate some of its basic properties. In particular we show that it combines well with the sum operation on graphs. These operations work on abstract graphs. Intuitively, concatenation is sequential composition of graphs, and sum is parallel composition of graphs.

If  $g$  and  $h$  are graphs with  $\text{type}(g) = (k, m)$  and  $\text{type}(h) = (m, n)$ , then their *concatenation*  $g \circ h$  is the graph obtained by first taking the disjoint union of  $g$  and  $h$ , and then identifying the  $i$ th end node of  $g$  with the  $i$ th begin node of  $h$ , for every  $i \in \{1, \dots, m\}$ ; moreover,  $\text{begin}(g \circ h) = \text{begin}(g)$  and  $\text{end}(g \circ h) = \text{end}(h)$ , and so  $\text{type}(g \circ h) = (k, n)$ . Note that the concatenation of  $g$  and  $h$  is defined only when  $|\text{end}(g)| = |\text{begin}(h)|$ . Formally, the definition is as follows (where we use node identification as defined in Section 2.2).

**Definition 3.** Let  $g$  and  $h$  be graphs such that  $|\text{end}(g)| = |\text{begin}(h)|$ . We assume that  $g$  and  $h$  are disjoint (otherwise an isomorphic copy of  $g$  or  $h$  should be taken). The *concatenation*  $g \circ h$  of  $g$  and  $h$  is the graph  $(g\&h)/R$  where  $g\&h = (V_g \cup V_h, E_g \cup E_h, s_g \cup s_h, t_g \cup t_h, l_g \cup l_h, \text{begin}(g), \text{end}(h))$  and  $R = \{(\text{end}(g)(i), \text{begin}(h)(i)) \mid 1 \leq i \leq |\text{end}(g)|\}$ .  $\square$

The *sum*  $g \oplus h$  of arbitrary graphs  $g$  and  $h$  is their disjoint union, with their sequences of begin nodes concatenated, and similarly for their end nodes. More formally, assuming that  $g$  and  $h$  are disjoint,  $g \oplus h = (V_g \cup V_h, E_g \cup E_h, s_g \cup s_h, t_g \cup t_h, l_g \cup l_h, \text{begin}(g) \cdot \text{begin}(h), \text{end}(g) \cdot \text{end}(h))$ .

The sum operation is taken from [BauCou] (where only graphs without end nodes are considered). All other operations in [BauCou] (viz. source redefinitions and source fusions) are unary operations, each of which is left-concatenation with a specific fixed graph.



**Fig. 1.** Two graphs, their concatenation, and their sum.

Figure 1 shows two (ordinary) abstract graphs,  $g$  of type  $(2, 3)$  and  $h$  of type  $(3, 1)$ , with their concatenation  $g \circ h$  of type  $(2, 1)$  and their sum  $g \oplus h$  of type



(5, 4). The graphs are drawn in the usual way; the  $i$ th begin node is indicated by  $b_i$ , and the  $i$ th end node by  $e_i$ .

These two graph operations have a number of simple properties. We stress again that the following lemmas are all about abstract graphs; in particular, the equality sign refers to the equality of abstract graphs (which is isomorphism of concrete graphs). First of all we show the basic fact that replacements are homomorphisms with respect to concatenation (just as string homomorphisms, of which they are a generalization) and with respect to sum.

**Lemma 4.** *Let  $\phi : \Sigma \rightarrow \text{GR}$  be a replacement, and let  $g, h \in \text{GR}(\Sigma)$ .*

- (1) *if  $|\text{end}(g)| = |\text{begin}(h)|$ , then  $\phi(g \circ h) = \phi(g) \circ \phi(h)$ , and*
- (2)  *$\phi(g \oplus h) = \phi(g) \oplus \phi(h)$ .*

*Proof.* (1) It is easy to verify this equality in the case that both  $g$  and  $h$  are atoms (for the definition of an atom, see Section 2.2). The general case is then proved as follows. Let  $\sigma$  and  $\tau$  be two symbols with the same type as  $g$  and  $h$ , respectively. Let  $\psi : \{\sigma, \tau\} \rightarrow \text{GR}(\Sigma)$  be the replacement with  $\psi(\sigma) = g$  and  $\psi(\tau) = h$ . Then, by the above special case,  $\psi(\text{atom}(\sigma) \circ \text{atom}(\tau)) = \psi(\text{atom}(\sigma)) \circ \psi(\text{atom}(\tau)) = \psi(\sigma) \circ \psi(\tau) = g \circ h$ . Hence  $\phi(g \circ h) = \phi(\psi(\text{atom}(\sigma) \circ \text{atom}(\tau)))$ . By Proposition 1,  $\phi \circ \psi$  is a replacement. Hence, again by the above special case,  $\phi(g \circ h) = \phi(\psi(\text{atom}(\sigma)) \circ \psi(\text{atom}(\tau))) = \phi(\psi(\sigma)) \circ \phi(\psi(\tau)) = \phi(g) \circ \phi(h)$ .

The proof of (2) is analogous.  $\square$

This lemma allows us to prove laws about  $\circ$  and  $\oplus$  by proving them for atoms only (as, in fact, we already did in the proof of Lemma 4). The next lemma summarizes the main basic properties of  $\circ$  and  $\oplus$ .

**Definition 5.** For every  $n \in \mathbf{N}$  the *identity*  $\text{id}_n$  of type  $(n, n)$  is the discrete graph with nodes  $x_1, \dots, x_n$  and  $\text{begin}(\text{id}_n) = \text{end}(\text{id}_n) = x_1 \cdots x_n$ . Thus,  $\text{id}_n$  is the (abstract) graph  $(\{x_1, \dots, x_n\}, \emptyset, \emptyset, \emptyset, \emptyset, \langle x_1, \dots, x_n \rangle, \langle x_1, \dots, x_n \rangle)$ . In particular,  $\text{id}_0$  is the empty graph.  $\square$

**Lemma 6.**

- (1) *Concatenation is associative, i.e., if  $|\text{end}(g_1)| = |\text{begin}(g_2)|$  and  $|\text{end}(g_2)| = |\text{begin}(g_3)|$ , then  $(g_1 \circ g_2) \circ g_3 = g_1 \circ (g_2 \circ g_3)$ .*
- (2) *The  $\text{id}_n$  are identities with respect to concatenation, i.e.,  $g \circ \text{id}_n = g$  and  $\text{id}_n \circ h = h$  for every  $g$  with  $|\text{end}(g)| = n$  and  $h$  with  $|\text{begin}(h)| = n$ .*
- (3) *Sum is associative with unity  $\text{id}_0$ , i.e.,  $(g_1 \oplus g_2) \oplus g_3 = g_1 \oplus (g_2 \oplus g_3)$  and  $g \oplus \text{id}_0 = \text{id}_0 \oplus g = g$ .*
- (4) *For every  $m, n \in \mathbf{N}$ ,  $\text{id}_{m+n} = \text{id}_m \oplus \text{id}_n$ .*
- (5) *Concatenation and sum satisfy the law of strict monoidality: if  $|\text{end}(g)| = |\text{begin}(g')|$  and  $|\text{end}(h)| = |\text{begin}(h')|$ , then*

$$(g \oplus h) \circ (g' \oplus h') = (g \circ g') \oplus (h \circ h').$$

*Proof.* (1) It is easy to verify that concatenation is associative for atoms. Now let  $\sigma_i$  be a symbol with the same type as  $g_i$ , and let  $\phi$  be the replacement with  $\phi(\sigma_i) = g_i$ . Then  $\phi((\text{atom}(\sigma_1) \circ \text{atom}(\sigma_2)) \circ \text{atom}(\sigma_3)) = \phi(\text{atom}(\sigma_1) \circ$

$(\text{atom}(\sigma_2) \circ \text{atom}(\sigma_3))$ , and, by Lemma 4,  $\phi((\text{atom}(\sigma_1) \circ \text{atom}(\sigma_2)) \circ \text{atom}(\sigma_3)) = (\phi(\text{atom}(\sigma_1)) \circ \phi(\text{atom}(\sigma_2))) \circ \phi(\text{atom}(\sigma_3)) = (g_1 \circ g_2) \circ g_3$  and  $\phi(\text{atom}(\sigma_1) \circ (\text{atom}(\sigma_2) \circ \text{atom}(\sigma_3))) = \phi(\text{atom}(\sigma_1)) \circ (\phi(\text{atom}(\sigma_2)) \circ \phi(\text{atom}(\sigma_3))) = g_1 \circ (g_2 \circ g_3)$ .

Properties (2), (3), and (5) can be shown in exactly the same way: by verifying them for atoms, and applying Lemma 4. Note that  $\phi(\text{id}_n) = \text{id}_n$  for every replacement  $\phi$ . Property (4) is obvious.  $\square$

Lemma 6 means that GR is a strict monoidal category (or x-category), see, e.g., [EhrKKK, Hot1, Ben]. The objects of this category are the natural numbers in  $\mathbf{N}$ , and each (abstract) graph of type  $(m, n)$  is a morphism from  $m$  to  $n$  in this category. Concatenation is the composition of morphisms (but is usually written  $h \circ g$  rather than  $g \circ h$ ), and the  $\text{id}_n$  are the identity morphisms. The set of objects and the set of morphisms form a monoid with respect to  $+$  and  $\oplus$ , respectively (where  $+$  is ordinary addition for natural numbers, with monoid identity 0).

We now show that all graphs can be built from a small number of elementary graphs with the operations of concatenation and sum.

For  $m, n \in \mathbf{N}$ , let  $I_{m,n}$  be the graph of type  $(m, n)$  with one node  $x$ , no edges,  $\text{begin}(I_{m,n}) = x^m = \langle x, \dots, x \rangle$  ( $m$  times), and  $\text{end}(I_{m,n}) = x^n = \langle x, \dots, x \rangle$  ( $n$  times). Note that  $I_{1,1} = \text{id}_1$ . Let  $\pi_{12}$  be the graph of type  $(2, 2)$  with two nodes  $x$  and  $y$ , no edges,  $\text{begin}(\pi_{12}) = xy$ , and  $\text{end}(\pi_{12}) = yx$ . For every typed alphabet  $\Sigma$  we define the set of *elementary* graphs over  $\Sigma$  by

$$\text{EL}(\Sigma) = \{\text{atom}(\sigma) \mid \sigma \in \Sigma\} \cup \{I_{0,1}, I_{1,0}, I_{1,2}, I_{2,1}, \pi_{12}, \text{id}_0\}.$$

**Theorem 7.** *For every typed alphabet  $\Sigma$ ,  $\text{GR}(\Sigma)$  is the smallest class of graphs containing  $\text{EL}(\Sigma)$  and closed under  $\circ$  and  $\oplus$ .*

*Proof.* We have to show that every graph in  $\text{GR}(\Sigma)$  can be written as an expression with the operators  $\circ$  and  $\oplus$ , and constants from  $\text{EL}(\Sigma)$ . We do this by reducing the problem to smaller and smaller sets of graphs. First we reduce it to the class of discrete graphs, i.e., graphs without edges.

Let  $g \in \text{GR}(\Sigma)$ , and let  $e$  be an edge of  $g$  with  $l_g(e) = \sigma$ . We will remove  $e$  from  $g$ , and express  $g$  in terms of the so obtained graph  $g'$  that has one edge less than  $g$  (and in terms of discrete graphs). By repeating this procedure, we can express  $g$  in terms of discrete graphs only. Let  $g' = (V_g, E_g - \{e\}, s, t, l, \text{begin}(g), \text{end}(g) \cdot s_g(e) \cdot t_g(e))$  where  $s, t, l$  are the restrictions of  $s_g, t_g, l_g$  to  $E_g - \{e\}$ , respectively. Thus,  $g'$  is obtained from  $g$  by removing  $e$ ; moreover, in order to be able to reconstruct  $g$  from  $g'$ , the sources and targets of  $e$  are turned into end nodes. It is now easy to verify that

$$g = g' \circ (\text{id}_n \oplus \text{backfold}(\text{atom}(\sigma)))$$

where  $n = |\text{end}(g)|$ , and the backfold operation is the one defined at the end of Section 2.2. Intuitively, the end nodes of  $\text{atom}(\sigma)$  are turned into begin nodes (by the backfold operation), and then they are glued to the new end nodes of  $g'$ . It is easy to prove that, for every graph  $h$ ,

$$\text{backfold}(h) = (h \oplus \text{id}_q) \circ \text{backfold}(\text{id}_q)$$

where  $q = |\text{end}(h)|$ . Consequently,

$$g = g' \circ (\text{id}_n \oplus (\text{atom}(\sigma) \oplus \text{id}_q) \circ \text{backfold}(\text{id}_q))$$

where  $q = |t(e)|$ . This shows that  $g$  can be expressed in terms of  $g'$  and discrete graphs.

It remains to find an expression for every discrete graph. To this aim we define the following special *permutation graphs*. Let  $k \geq 1$  and let  $\alpha$  be a permutation of  $\{1, \dots, k\}$ . Then  $\pi_\alpha$  is the discrete graph with nodes  $\{x_1, \dots, x_k\}$ ,  $\text{begin}(\pi_\alpha) = x_1 \cdots x_k$ , and  $\text{end}(\pi_\alpha) = x_{\alpha(1)} \cdots x_{\alpha(k)}$ . Note that  $\pi_{12}$  is the permutation graph  $\pi_\alpha$  with  $\alpha(1) = 2$  and  $\alpha(2) = 1$ . We need some simple properties of permutation graphs. In what follows we write  $[n]$  for  $\{1, \dots, n\}$ , for every  $n \in \mathbf{N}$ . First, if  $\alpha$  and  $\beta$  are permutations of  $[k]$ , then  $\pi_\alpha \circ \pi_\beta = \pi_{\alpha \circ \beta}$ , and if  $\text{id}$  is the identity permutation of  $[k]$ , then  $\pi_{\text{id}} = \text{id}_k$ . Second, let  $g$  be a graph of type  $(m, n)$  with  $V_g = \{x_1, \dots, x_k\}$ ,  $\text{begin}(g) = x_{\gamma(1)} \cdots x_{\gamma(m)}$ , and  $\text{end}(g) = x_{\delta(1)} \cdots x_{\delta(n)}$ , where  $\gamma : [m] \rightarrow [k]$  and  $\delta : [n] \rightarrow [k]$ . If  $\alpha$  is a permutation of  $[n]$ , then  $g \circ \pi_\alpha$  is the same graph as  $g$  except that  $\text{end}(g \circ \pi_\alpha) = x_{\delta(\alpha(1))} \cdots x_{\delta(\alpha(n))}$ . This means that  $g \circ \pi_\alpha$  is obtained from  $g$  by applying permutation  $\alpha$  to  $\text{end}(g)$ . Similarly, if  $\alpha$  is a permutation of  $[m]$ , then  $\pi_\alpha \circ g$  is the same graph as  $g$  except that  $\text{begin}(\pi_\alpha \circ g) = x_{\gamma(\alpha^{-1}(1))} \cdots x_{\gamma(\alpha^{-1}(m))}$ . Thus, to obtain  $\pi_\alpha \circ g$  from  $g$ , permutation  $\alpha^{-1}$  is applied to  $\text{begin}(g)$ .

Now let  $g$  be an arbitrary discrete graph, with  $\text{type}(g) = (m, n)$ ,  $V_g = \{x_1, \dots, x_k\}$ ,  $\text{begin}(g) = x_{\gamma(1)} \cdots x_{\gamma(m)}$ , and  $\text{end}(g) = x_{\delta(1)} \cdots x_{\delta(n)}$ , where  $\gamma : [m] \rightarrow [k]$  and  $\delta : [n] \rightarrow [k]$ . For every  $1 \leq i \leq k$ , let  $p_i$  be the number of occurrences of  $x_i$  in  $\text{begin}(g)$ , and let  $q_i$  be the number of occurrences of  $x_i$  in  $\text{end}(g)$ . Let  $\alpha$  be any permutation of  $[m]$  such that  $x_{\gamma(\alpha(1))} \cdots x_{\gamma(\alpha(m))} = x_1^{p_1} \cdots x_k^{p_k}$ , and let  $\beta$  be any permutation of  $[n]$  such that  $x_{\delta(\beta(1))} \cdots x_{\delta(\beta(n))} = x_1^{q_1} \cdots x_k^{q_k}$ . Thus, intuitively,  $\alpha$  and  $\beta$  order  $\text{begin}(g)$  and  $\text{end}(g)$ , respectively. Clearly, by the above properties of permutation graphs, the graph  $\pi_{\alpha^{-1}} \circ g \circ \pi_\beta$  has the same nodes as  $g$ , has begin nodes  $x_1^{p_1} \cdots x_k^{p_k}$ , and has end nodes  $x_1^{q_1} \cdots x_k^{q_k}$ . Hence  $\pi_{\alpha^{-1}} \circ g \circ \pi_\beta = I_{p_1, q_1} \oplus \cdots \oplus I_{p_k, q_k}$ . By multiplying with  $\pi_\alpha$  to the left, and with  $\pi_{\beta^{-1}}$  to the right, we obtain that  $g = \pi_\alpha \circ (I_{p_1, q_1} \oplus \cdots \oplus I_{p_k, q_k}) \circ \pi_{\beta^{-1}}$ .

It now remains to find expressions for all graphs  $I_{m,n}$  and all graphs  $\pi_\alpha$ . The following equations show how to find an expression for  $I_{m,n}$ :

$$\begin{aligned} I_{1,1} &= I_{1,2} \circ I_{2,1} \\ I_{m+1,1} &= (I_{m,1} \oplus I_{1,1}) \circ I_{2,1} \quad \text{for } m \geq 2 \\ I_{1,n+1} &= I_{1,2} \circ (I_{1,n} \oplus I_{1,1}) \quad \text{for } n \geq 2 \\ I_{m,n} &= I_{m,1} \circ I_{1,n} \quad \text{for } m, n \in \mathbf{N}. \end{aligned}$$

Clearly, the identity graphs can also be expressed: for every  $n \in \mathbf{N}$ ,  $\text{id}_n = I_{1,1} \oplus \cdots \oplus I_{1,1}$  ( $n$  times). To find an expression for  $\pi_\alpha$ , where  $\alpha$  is a permutation of  $[k]$ , we note that either  $\alpha$  is the identity on  $[k]$ , in which case  $\pi_\alpha = \text{id}_k$ , or  $\alpha$  is the composition of interchanging permutations, where an interchanging permutation  $\alpha_i$  interchanges  $i$  and  $i+1$  and leaves the other numbers as they are (with  $1 \leq i \leq k-1$ ). In the latter case, by the property of permutation

graphs mentioned above,  $\pi_\alpha$  is the concatenation of graphs  $\pi_{\alpha_i}$ . Now, clearly,  $\pi_{\alpha_i} = \text{id}_{i-1} \oplus \pi_{12} \oplus \text{id}_{k-i-1}$ .

This shows that all graphs in  $\text{GR}(\Sigma)$  can be expressed in terms of  $\circ$ ,  $\oplus$ , and the constants in  $\text{EL}(\Sigma)$ .  $\square$

Theorem 7 is analogous to Proposition 3.6 of [BauCou]. It is open whether there exists a complete set of equations (including those of Lemma 6) for the operations in  $\{\circ, \oplus\} \cup \text{EL}(\Sigma)$ . This would give a result analogous to Theorem 3.10 of [BauCou]. It would characterize  $\text{GR}(\Sigma)$  as the free x-category satisfying the equations; such results are shown in [Hot1, Cla] (where  $I_{1,0}, I_{1,2}, \pi_{12}$  are denoted  $U, D, V$ , respectively).

It is not difficult to show that the set  $\text{EL}(\Sigma)$  is minimal, in the sense that if one removes one element from it, then Theorem 7 does not hold any more. Also, it should be clear that the concatenation operation cannot be dropped from Theorem 7, even if one would replace  $\text{EL}(\Sigma)$  by another finite set of graphs (because, with sum, only graphs with very small connected components could be built). To show that the sum operation cannot be dropped from Theorem 7, we now discuss the close relationship between the concatenation operation and the notion of pathwidth (introduced in [RobSey]; see also, e.g., [Bod, Klo, EllST]). In the following definition we (slightly) generalize the notion of pathwidth, to (hyper)graphs with begin and end nodes (cf. [Cou3]).

**Definition 8.** A *path decomposition* of a graph  $g$  is a sequence  $(V_1, \dots, V_n)$ ,  $n \geq 1$ , of subsets of  $V_g$  such that

- (1)  $\bigcup_{i=1}^n V_i = V_g$ ,
- (2) for every  $e \in E_g$  there is an  $i$  with  $s(e) \in V_i^*$  and  $t(e) \in V_i^*$ ,
- (3) if  $i < k < j$ , then  $V_i \cap V_j \subseteq V_k$ , and
- (4)  $\text{begin}(g) \in V_1^*$  and  $\text{end}(g) \in V_n^*$ .

The *width* of  $(V_1, \dots, V_n)$  is  $\max\{\#V_i \mid 1 \leq i \leq n\} - 1$ , where  $\#V_i$  is the cardinality of  $V_i$ .

The *pathwidth* of a graph  $g$ , denoted  $\text{pathwidth}(g)$ , is the minimal width of a path decomposition of  $g$ .  $\square$

The relationship between concatenation and pathwidth is expressed in the following result, which (in view of Theorem 23) is essentially due to [Lau] (see also [Cou3]).

**Theorem 9.** *Let  $k \geq 1$ . For every graph  $g$ ,  $\text{pathwidth}(g) \leq k$  if and only if there exist graphs  $g_1, \dots, g_n$ ,  $n \geq 1$ , such that  $g = g_1 \circ \dots \circ g_n$  and  $\#V_{g_i} \leq k + 1$  for every  $1 \leq i \leq n$ .*

*Proof.* We prove by induction on  $n$  that  $g$  has a path decomposition  $(V_1, \dots, V_n)$  of width  $\leq k$  if and only if there exist graphs  $g_1, \dots, g_n$  with  $\#V_{g_i} \leq k + 1$  such that  $g = g_1 \circ \dots \circ g_n$ . For  $n = 1$  this is obvious.

Assume that  $g$  has a path decomposition  $(V_1, \dots, V_n, V_{n+1})$  with  $\#V_i \leq k + 1$ . By condition (3) of Definition 8,  $(V_1 \cup \dots \cup V_n) \cap V_{n+1} = V_n \cap V_{n+1}$ . Let  $g'$  be the subgraph of  $g$  induced by  $V_1 \cup \dots \cup V_n$ , such that  $\text{begin}(g') = \text{begin}(g)$  and  $\text{end}(g)$

consists of the nodes in  $V_n \cap V_{n+1}$ , in some order. Let  $g_{n+1}$  be the subgraph of  $g$  induced by  $V_{n+1}$  with  $\text{begin}(g_{n+1}) = \text{end}(g')$  and  $\text{end}(g_{n+1}) = \text{end}(g)$ . Clearly,  $g = g' \circ g_{n+1}$ . Also,  $(V_1, \dots, V_n)$  is a path decomposition of  $g'$ , and hence, by induction  $g' = g_1 \circ \dots \circ g_n$  with  $\#V_{g_i} \leq k + 1$ , and so  $g = g_1 \circ \dots \circ g_n \circ g_{n+1}$ .

Assume now that  $g = g_1 \circ \dots \circ g_n \circ g_{n+1}$  with  $\#V_{g_i} \leq k + 1$ . Let  $g' = g_1 \circ \dots \circ g_n$ . Hence  $g = g' \circ g_{n+1}$ . We may assume that  $g'$  and  $g_{n+1}$  are disjoint, and that  $g = (g' \& g_{n+1})/R$ , as in Definition 3. By induction,  $g'$  has a path decomposition  $(V_1, \dots, V_n)$  with  $\#V_i \leq k + 1$ . Let  $V_{n+1} = V_{g_{n+1}}$ . It should now be clear that the sequence  $(V_1/R, \dots, V_n/R, V_{n+1}/R)$  is a path decomposition of  $g$ .  $\square$

Since there are graphs of arbitrary large pathwidth (such as the complete graph on  $n$  nodes, which has pathwidth  $n - 1$ ), this theorem implies that, for any typed alphabet  $\Sigma$ , there is no finite subset  $E$  of  $\text{GR}(\Sigma)$  such that  $\text{GR}(\Sigma)$  is the smallest set of graphs containing  $E$  and closed under concatenation (because the pathwidth of all graphs in this smallest set is at most equal to the maximal size of the graphs in  $E$ ).

## 4 Context-free graph grammars

In this section we use context-free grammars to generate graph expressions that are built from arbitrary constant graphs with the graph operators  $\circ$  and  $\oplus$ . Taking the values of these expressions in  $\text{GR}$ , each such context-free grammar generates a graph language.

Let  $\text{CS}$  be the set of operators  $\{\circ, \oplus\} \cup \{c_g \mid g \in \text{GR}\}$ , where  $\circ$  and  $\oplus$  denote concatenation and sum of graphs, as usual, and  $c_g$  is a constant standing for the graph  $g$ .

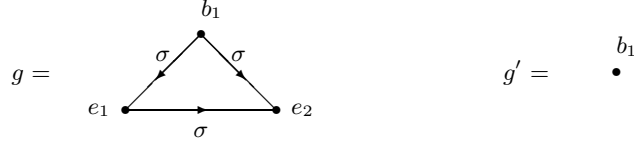
Expressions over  $\text{CS}$  are defined in the usual way. Let  $\Sigma$  be a typed alphabet, disjoint with  $\text{CS}$  (where, intuitively, each  $\sigma \in \Sigma$  is a variable that ranges over all graphs with the same type as  $\sigma$ ). A (well-formed) *expression* over  $\text{CS}$  and  $\Sigma$  is a string over  $\text{CS} \cup \Sigma \cup \{(\cdot, \cdot)\}$  defined recursively as follows, together with its *type*: (1) every  $\sigma \in \Sigma$  is an expression, with the same type, (2) every constant  $c_g$  is an expression, with  $\text{type}(c_g) = \text{type}(g)$ , (3) if  $e$  and  $f$  are expressions with  $\text{type}(e) = (k, m)$  and  $\text{type}(f) = (m, n)$ , then  $(e \circ f)$  is an expression with  $\text{type}(e \circ f) = (k, n)$ , and (3) if  $e$  and  $f$  are expressions with  $\text{type}(e) = (m, n)$  and  $\text{type}(f) = (p, q)$ , then  $(e \oplus f)$  is an expression with  $\text{type}(e \oplus f) = (m + p, n + q)$ . An ‘expression over  $\text{CS}$ ’ is defined in the same way, without clause (1). If  $e$  is an expression over  $\text{CS}$ , then its *value*, denoted by  $\text{val}(e)$ , is a graph in  $\text{GR}$ , defined recursively in the usual way:  $\text{val}(c_g) = g$ ,  $\text{val}(e \circ f) = \text{val}(e) \circ \text{val}(f)$ , and  $\text{val}(e \oplus f) = \text{val}(e) \oplus \text{val}(f)$ .

**Definition 10.** A *context-free graph grammar* over  $\text{CS}$  is an ordinary context-free grammar  $G = (N, T, P, S)$ , see Section 2.1, such that  $N$  is a typed alphabet,  $T$  is a finite subset of  $\text{CS} \cup \{(\cdot, \cdot)\}$ , and the right-hand side of each production in  $P$  is an expression over  $\text{CS}$  and  $N$ , of the same type as the left-hand side.  $\square$

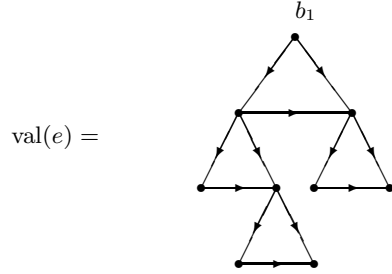
Obviously, the context-free language  $L(G)$  generated by  $G$  is a set of expressions over  $\text{CS}$  (and it is also a regular tree language, see [GécSte]). The *graph language*

generated by  $G$  is  $\text{val}(L(G)) = \{\text{val}(e) \mid e \in L(G)\}$ . Note that  $\text{type}(\text{val}(L(G))) = \text{type}(S)$ . By  $\text{Val}(\text{CFG}(\text{CS}))$  we denote the class of all graph languages generated by context-free graph grammars over CS. By the results of [MezWri], it is the class of *equational* subsets of the algebra of graphs with the operations  $\circ$  and  $\oplus$ .

It should be clear that, due to Theorem 7, we could restrict CS to contain only elementary constants, i.e., constants  $c_g$  with  $g \in \text{EL}(\Sigma)$  for some  $\Sigma$ . This would give the same class  $\text{Val}(\text{CFG}(\text{CS}))$ .



**Fig. 2.** Graphs  $g$  and  $g'$ .



**Fig. 3.** The value of graph expression  $e$ .

As an example, consider the context-free graph grammar  $G_b$  that has one nonterminal  $X$ , with  $\text{type}(X) = (1, 0)$ , and two productions  $X \rightarrow c_g \circ (X \oplus X)$  and  $X \rightarrow c_{g'}$ , where  $g$  is the triangle of type  $(1, 2)$  with  $V = \{x, y, z\}$ ,  $E = \{(x, y), (x, z), (y, z)\}$ ,  $s(u, v) = u$ ,  $t(u, v) = v$ , and  $l(u, v) = \sigma$  for every edge  $(u, v)$ ,  $\text{begin}(g) = x$  and  $\text{end}(g) = yz$ , and  $g'$  is the graph of type  $(1, 0)$  with one node  $x$ , no edges,  $\text{begin}(g') = x$  and  $\text{end}(g') = \lambda$ . The graphs  $g$  and  $g'$  are shown in Fig. 2. The expression

$$e = c_g \circ (c_g \circ (c_{g'} \oplus c_g \circ (c_{g'} \oplus c_{g'})) \oplus c_g \circ (c_{g'} \oplus c_{g'}))$$

is in  $L(G)$ ; the graph  $\text{val}(e)$  is shown in Fig. 3 (without the edge labels  $\sigma$ ). Clearly,  $\text{val}(L(G_b))$  is the set of all graphs of type  $(1, 0)$  that are obtained from (directed, rooted) binary trees by connecting each pair of children by an additional edge; the sequence of begin nodes consists of the root of the binary tree. This graph language is therefore in  $\text{Val}(\text{CFG}(\text{CS}))$ .

The main result of this section is that generating graph languages in the above way is equivalent to generating them with HR grammars (see Section 2.3).

Thus, the HR grammars generate exactly the equational subsets of the algebra of graphs with the operations  $\circ$  and  $\oplus$ . As observed in the introduction, this is a simple variant of Theorem 4.11 of [BauCou] (and the proof is analogous).

**Theorem 11.**  $\text{Val}(\text{CFG}(\text{CS})) = \text{HR}$ .

*Proof.* Similar to the restriction on productions of HR grammars, we can assume without loss of generality that no nonterminal occurs more than once in the right-hand side of a production of a context-free graph grammar. Moreover, we can also assume that, in a context-free graph grammar, the nonterminals do not occur as edge labels in the constants  $c_g$  that are used in the right-hand sides of its productions.

To turn a context-free graph grammar into an HR grammar, we extend the definition of the ‘val’ function to expressions over CS and  $N$ . This is simply done by extending the recursive definition of ‘val’ with the requirement that  $\text{val}(X) = \text{atom}(X)$  for every  $X \in N$ .

Let  $G$  be a context-free graph grammar and  $G'$  an HR grammar. We say that  $G$  and  $G'$  are *related* if they have the same typed alphabet of nonterminals, with the same initial nonterminal, and  $P' = \{X \rightarrow \text{val}(t) \mid X \rightarrow t \in P\}$ , where  $P$  is the set of productions of  $G$  and  $P'$  the one of  $G'$ . Trivially, for every context-free graph grammar there is a related HR grammar. The other way around, it suffices to show that for every graph  $h \in \text{GR}(N \cup T)$ , where  $N$  and  $T$  are the terminal and nonterminal alphabet of the HR grammar, respectively, there is an expression  $t$  over CS and  $N$  such that  $\text{val}(t) = h$ . By Theorem 7 there is an expression  $e$  over CS such that  $\text{val}(e) = h$ , and for every constant  $c_g$  that occurs in  $e$ ,  $g \in \text{EL}(N \cup T)$ . Let  $t$  be the expression that is obtained from  $e$  by changing every subexpression  $\text{atom}(X)$  into  $X$ , for every  $X \in N$ . Obviously  $t$  is the required expression. Hence, for every HR grammar there is a related context-free graph grammar.

It now suffices to show that related grammars  $G$  and  $G'$  generate the same graph language. To this aim we show that for every nonterminal  $X$  and every terminal graph  $g$ ,  $\text{atom}(X) \Rightarrow^* g$  in  $G'$  if and only if there exists an expression  $e$  over CS such that  $X \Rightarrow^* e$  in  $G$  and  $\text{val}(e) = g$ . This can be proved by induction on the length of the derivations, as follows.

Consider a derivation  $X \Rightarrow t \Rightarrow^* e$  in  $G$ . Let  $t$  contain the nonterminals  $X_1, \dots, X_n$  (and recall that each nonterminal  $X_i$  occurs exactly once in  $t$ ). Then there exist expressions  $e_i$  such that  $X_i \Rightarrow^* e_i$  and  $e = \psi(t)$  where  $\psi$  is the string homomorphism with  $\psi(X_i) = e_i$  and the identity otherwise. Now let  $\phi$  be the replacement with  $\phi(X_i) = \text{val}(e_i)$  and  $\phi(\sigma) = \text{atom}(\sigma)$  for all terminal symbols.

It is straightforward to show that  $\text{val}(\psi(t)) = \phi(\text{val}(t))$ , by induction on the structure of the expression  $t$  (cf. Proposition 4.7 of [BauCou]). As an example, if  $t = t_1 \circ t_2$ , then  $\text{val}(\psi(t)) = \text{val}(\psi(t_1) \circ \psi(t_2)) = \text{val}(\psi(t_1)) \circ \text{val}(\psi(t_2)) = \phi(\text{val}(t_1)) \circ \phi(\text{val}(t_2)) = \phi(\text{val}(t_1) \circ \text{val}(t_2)) = \phi(\text{val}(t))$ , where we have used Lemma 4(1). As another example, if  $t = X_i$ , then  $\phi(\text{val}(t)) = \phi(\text{atom}(X_i)) = \phi(X_i) = \text{val}(e_i) = \text{val}(\psi(t))$ .

By induction,  $\text{atom}(X_i) \Rightarrow^* \text{val}(e_i)$  in  $G'$ . Since  $G$  and  $G'$  are related,  $G'$  has the production  $X \rightarrow \text{val}(t)$ . It is easy to see that  $\text{val}(t)$  has  $n$  nonterminal edges,

labeled by  $X_1, \dots, X_n$ . Hence, by Proposition 2,  $\text{atom}(X) \Rightarrow \text{val}(t) \Rightarrow^* \phi(\text{val}(t))$ . This shows that  $\text{atom}(X) \Rightarrow^* \text{val}(\psi(t)) = \text{val}(e)$ .

The proof in the other direction is similar and is left to the reader.  $\square$

Since the concept of *related* grammars, as discussed in the above proof, preserves the number of nonterminals in the right-hand sides of productions, the above result is also true in the linear case. By  $\text{Val}(\text{LIN-CFG}(\text{CS}))$  we denote the class of languages generated by linear context-free graph grammars over CS.

**Corollary 12.**  $\text{Val}(\text{LIN-CFG}(\text{CS})) = \text{LIN-HR}$ .

However, in the linear case, the form of the context-free graph grammar can even be restricted to be “right-linear” in the following sense. A context-free graph grammar over CS is *right-linear* if its productions are of the form  $X \rightarrow c_g \circ Y$  or of the form  $X \rightarrow c_g$ , where  $X$  and  $Y$  are nonterminals. Note in particular that  $\oplus$  is not needed. By  $\text{Val}(\text{RLIN-CFG}(\text{CS}))$  we denote the class of graph languages generated by right-linear context-free graph grammars over CS.

**Theorem 13.**  $\text{Val}(\text{RLIN-CFG}(\text{CS})) = \text{LIN-HR}$ .

*Proof.* By Corollary 12, it suffices to show that  $\text{LIN-HR} \subseteq \text{Val}(\text{RLIN-CFG}(\text{CS}))$ . Let  $L$  be a graph language in LIN-HR, and let  $G = (N, T, P, S)$  be a linear HR grammar generating  $L$ .

We first consider the case that for every  $X \in N$  there exists  $m \in \mathbf{N}$  such that  $\text{type}(X) = (m, 0)$ . By the proof of Theorem 11 it suffices to construct a context-free graph grammar  $G'$  that is related to  $G$ .  $G'$  has the same nonterminal alphabet as  $G$ , with the same initial nonterminal.  $G'$  has the set of productions  $P' = \{p' \mid p \in P\}$ , where, for each  $p \in P$ ,  $p'$  is defined as follows. Let  $p$  be the production  $X \rightarrow g$ . If  $g \in \text{GR}(T)$ , then we define  $p'$  to be  $X \rightarrow c_g$ . Otherwise,  $g$  has exactly one edge  $e$  that is labeled with a nonterminal, say,  $Y$ . Note that  $\text{end}(g) = \lambda$  and  $t_g(e) = \lambda$ . Then we define  $p'$  to be the production  $X \rightarrow c_{g'} \circ Y$ , where  $g' = (V_g, E_g - \{e\}, s, t, l, \text{begin}(g), s_g(e))$  and  $s, t, l$  are the restrictions of  $s_g, t_g, l_g$  to  $E_g - \{e\}$ . Clearly,  $\text{val}(c_{g'} \circ Y) = g' \circ \text{atom}(Y) = g$ . Hence  $G$  and  $G'$  are indeed related. Note that the construction of  $g'$  from  $g$  is a special case of that in the first part of the proof of Theorem 7.

We now consider the general case. To be able to use the above special case, define the LIN-HR grammar  $\overline{G} = (\overline{N}, T, \overline{P}, S)$ , where  $\overline{N}$  is the same set as  $N$  with a different type function: if  $\text{type}(X) = (m, n)$  in  $N$ , then  $\text{type}(X) = (m + n, 0)$  in  $\overline{N}$ . For every graph  $h \in \text{GR}(N \cup T)$  we define the graph  $\overline{h} = (V_h, E_h, s, t, l_h, \text{begin}(h) \cdot \text{end}(h), \lambda)$  where, for  $e \in E_h$ ,  $s(e)$  and  $t(e)$  are defined as follows: if  $l_h(e) \in T$ , then  $s(e) = s_h(e)$  and  $t(e) = t_h(e)$ ; if  $l_h(e) \in N$ , then  $s(e) = s_h(e) \cdot t_h(e)$  and  $t(e) = \lambda$ . Note that if  $h \in \text{GR}(T)$  then  $\overline{h} = \text{backfold}(h)$ . We now define  $\overline{P} = \{X \rightarrow \overline{h} \mid X \rightarrow h \in P\}$ . This ends the definition of  $\overline{G}$ . It is straightforward to show that  $L(\overline{G}) = \{\text{backfold}(g) \mid g \in L(G)\}$ . In fact, the derivations of  $\overline{G}$  are exactly all  $\text{atom}(S) \Rightarrow \overline{g}_1 \Rightarrow \overline{g}_2 \Rightarrow \dots \Rightarrow \overline{g}_n$  where  $\text{atom}(S) \Rightarrow g_1 \Rightarrow g_2 \Rightarrow \dots \Rightarrow g_n$  is a derivation of  $G$ . Since  $\overline{G}$  satisfies the above special case, we conclude that  $\text{backfold}(L)$  is in  $\text{Val}(\text{RLIN-CFG}(\text{CS}))$ .



Suppose that  $\text{type}(L) = (m, n)$ . It is easy to verify, for every  $g$  of type  $(m, n)$ , that

$$g = (\text{id}_m \oplus \text{fold}(\text{id}_n)) \circ (\text{backfold}(g) \oplus \text{id}_n).$$

Hence  $L = \{(\text{id}_m \oplus \text{fold}(\text{id}_n)) \circ (h \oplus \text{id}_n) \mid h \in \text{backfold}(L)\}$ . Thus, it now suffices to show that if  $L'$  is in  $\text{Val}(\text{RLIN-CFG}(\text{CS}))$ , then so are all languages  $\{h \oplus \text{id}_n \mid h \in L'\}$ , for  $n \in \mathbf{N}$ , and all languages  $\{g_0 \circ h \mid h \in L'\}$ , for  $g_0 \in \text{GR}$ . To this aim, let  $G'$  be a right-linear context-free graph grammar generating  $L'$ .

Change, in the productions of  $G'$ , every constant  $c_g$  into the constant  $c_{g \oplus \text{id}_n}$ . Clearly, the resulting right-linear context-free graph grammar generates all graphs  $(g_1 \oplus \text{id}_n) \circ \dots \circ (g_k \oplus \text{id}_n)$  with  $g_1 \circ \dots \circ g_k \in L'$ . By the law of strict monoidality (Lemma 6(5)),

$$(g_1 \oplus \text{id}_n) \circ \dots \circ (g_k \oplus \text{id}_n) = (g_1 \circ \dots \circ g_k) \oplus (\text{id}_n \circ \dots \circ \text{id}_n) = (g_1 \circ \dots \circ g_k) \oplus \text{id}_n.$$

This proves that the resulting grammar generates  $\{g \oplus \text{id}_n \mid g \in L'\}$ .

Introduce a new initial nonterminal  $S'$ , and add to  $G'$  all the productions  $S' \rightarrow c_{g_0 \circ g} \circ Y$  and  $S' \rightarrow c_{g_0 \circ g}$  such that  $S \rightarrow c_g \circ Y$  and  $S \rightarrow c_g$  are productions of  $G'$ , respectively (where  $S$  is the old initial nonterminal of  $G'$ ). Clearly, the resulting right-linear grammar generates all graphs  $(g_0 \circ g_1) \circ g_2 \circ \dots \circ g_k$  with  $g_1 \circ \dots \circ g_k \in L'$ . In other words (using the associativity of concatenation), it generates the graph language  $\{g_0 \circ h \mid h \in L'\}$ . Note that we could also have added the one production  $S' \rightarrow c_{g_0} \circ S$ ; the reason for not doing so will become clear in the proof of Theorem 27.  $\square$

This result suggests that for context-free graph grammars there is no difference between the linear and the right-linear case, as opposed to the case of ordinary context-free grammars (where the right-linear grammars generate the regular languages which form a proper subclass of the linear languages). More support for this intuition will be given in the next section.

## 5 Strings Denote Graphs

Since concatenation of graphs is associative, strings can be viewed as expressions that denote graphs. Thus, as an even simpler variation of the approach with expression generating context-free grammars in Section 4, we can use all possible string grammars to generate graph languages. More generally, every class  $K$  of string languages defines a class  $\text{Int}(K)$  of graph languages (where  $\text{Int}$  stands for ‘interpretation’, which is similar to  $\text{Val}$  in Section 4). An ‘interpretation’ is a mapping that associates a graph with each symbol of an alphabet.

**Definition 14.** Let  $A$  be an alphabet. An *interpretation* of  $A$  is a mapping  $h : A \rightarrow \text{GR}$ ;  $h$  is extended to a (partial) function from  $A^*$  to  $\text{GR}$  by

$$h(a_1 a_2 \dots a_n) = h(a_1) \circ h(a_2) \circ \dots \circ h(a_n)$$

with  $n \geq 1$  and  $a_i \in A$  for all  $1 \leq i \leq n$ .  $\square$

Note that the extended  $h$  is partial because the types of the  $h(a_i)$  may not fit; moreover,  $h(\lambda)$  is undefined (where  $\lambda$  is the empty string). Thus, the only “technical trouble” is that the concatenation of graphs is typed whereas the concatenation of strings is always possible. To deal with this, the following lemma is useful. It says that the domain of an interpretation is regular.

**Lemma 15.** *For every interpretation  $h : A \rightarrow \text{GR}$ , the language  $\{w \in A^* \mid h(w) \text{ is defined}\}$  is regular.*

*Proof.* Clearly,  $h(a_1 a_2 \cdots a_n)$  is defined if and only if  $n \geq 1$  and  $|\text{end}(h(a_i))| = |\text{begin}(h(a_{i+1}))|$  for every  $1 \leq i < n$ . It is easy to construct a finite automaton that checks this.  $\square$

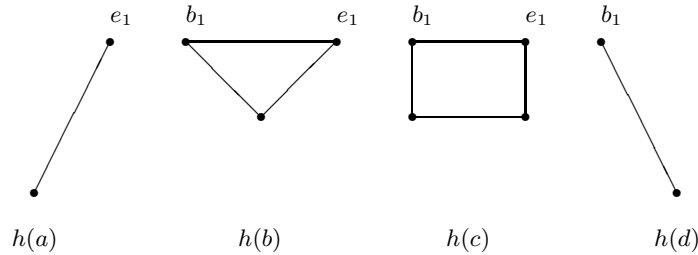
For a string language  $L \subseteq A^*$ , we define, as usual, the set of graphs  $h(L) = \{g \in \text{GR} \mid g = h(w) \text{ for some } w \in L\}$ ; note that  $h(L)$  need not be a graph language (in our particular meaning of the term, as defined in Section 2.2) because not all graphs need have the same type.

**Definition 16.** Let  $K$  be a class of string languages. The *interpretation* of  $K$  is  $\text{Int}(K) = \{h(L) \mid L \in K, h : A \rightarrow \text{GR} \text{ with } L \subseteq A^*, h(L) \text{ is a graph language}\}$ .  $\square$

In other words,  $\text{Int}(K)$  consists of all graph languages  $h(L)$ , where  $L$  is any language in  $K$  and  $h$  is any mapping from the symbols of  $L$  to graphs. Intuitively,  $h$  determines the interpretation of the symbols, and then the concatenation of those symbols is interpreted as concatenation of the corresponding graphs.

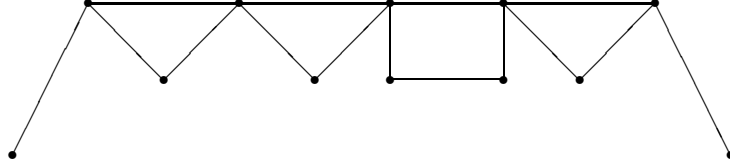
It is an immediate consequence of Theorem 9 that every graph language  $h(L)$  in  $\text{Int}(K)$  is of *bounded pathwidth*, i.e., there exists  $k$  such that  $\text{pathwidth}(g) \leq k$  for every  $g \in h(L)$ . In fact, if  $L \subseteq A^*$ , then  $k = \max\{\#V_{h(a)} \mid a \in A\} - 1$ .

**Corollary 17.** *For every  $K$ , every graph language in  $\text{Int}(K)$  is of bounded pathwidth.*



**Fig. 4.** An interpretation.

The first class  $K$  of interest is the class REG of regular languages. An example of a graph language in  $\text{Int}(\text{REG})$ , of type  $(0, 0)$ , is  $h(a(b \cup c)^*d)$  where the graphs  $h(a)$ ,  $h(b)$ ,  $h(c)$ , and  $h(d)$  are shown in Fig. 4 (without edge directions and edge labels). The graph  $h(abbcbd)$  is shown in Fig. 5. Clearly, the graph language  $h(a(b \cup c)^*d)$  consists of all “clothes lines” on which triangles and rectangles are hanging to dry. We first present a characterization of  $\text{Int}(\text{REG})$  by regular ex-



**Fig. 5.** Graph interpretation of the string  $abbcbd$ .

pressions, corresponding to the characterization of REG by regular expressions. To this aim we define the operations of union, concatenation, and (Kleene) star for graph languages. The operation of graph concatenation is extended to graph languages  $L$  and  $L'$  in the usual way: if  $\text{type}(L) = (k, m)$  and  $\text{type}(L') = (m, n)$ , then their *concatenation* is defined by  $L \circ L' = \{g \circ g' \mid g \in L, g' \in L'\}$ . Then, in the obvious way, the *star* of a graph language is defined by iterated concatenation: for a graph language  $L$  with  $\text{type}(L) = (k, k)$  for some  $k \in \mathbf{N}$ ,  $L^* = \bigcup_{n \in \mathbf{N}} L^n$  where  $L^n = L \circ \dots \circ L$  ( $n$  times) for  $n \geq 1$ , and  $L^0 = \{\text{id}_k\}$ . Also,  $L^+ = \bigcup_{n \geq 1} L^n$  is the (Kleene) plus of  $L$ . Finally, the *union*  $L \cup L'$  of two graph languages  $L$  and  $L'$  is defined only when  $\text{type}(L) = \text{type}(L')$  (otherwise it would not be a graph language). Thus, the operations of union, concatenation, and star are also typed operations on graph languages (as opposed to the case of string languages for which they are always defined). Let  $\text{REX}(\cup, \circ, *, \text{SING})$  denote the smallest class of graph languages containing the empty graph language and all singleton graph languages, and closed under the operations union, concatenation, and star. Thus, it is the class of all graph languages that can be denoted by (the usual) regular expressions, where the symbols of the alphabet denote singleton graph languages. As an example, the above graph language of clothes lines is in  $\text{REX}(\cup, \circ, *, \text{SING})$  because it can be written as  $\{h(a)\} \circ (\{h(b)\} \cup \{h(c)\})^* \circ \{h(d)\}$ .

**Theorem 18.**  $\text{Int}(\text{REG}) = \text{REX}(\cup, \circ, *, \text{SING})$ .

*Proof.* We have to cope with the “technical trouble” of typing, in particular with the empty string. Note that, for a graph language  $L$  with  $\text{type}(L) = (k, k)$ ,  $L^* = L^+ \cup \{\text{id}_k\}$  and  $L^+ = L \circ L^*$ . This shows that we can replace star by plus, i.e.,  $\text{REX}(\cup, \circ, *, \text{SING}) = \text{REX}(\cup, \circ, +, \text{SING})$ , the smallest class of graph languages

containing the empty graph language and all singleton graph languages, and closed under the operations union, concatenation, and plus.

To show that  $\text{REX}(\cup, \circ, +, \text{SING}) \subseteq \text{Int}(\text{REG})$ , it suffices to prove that  $\text{Int}(\text{REG})$  contains the empty language and all singleton graph languages, and that it is closed under union, concatenation, and plus. Clearly,  $h(\emptyset) = \emptyset$  for any interpretation  $h$ . Also, if  $h(a) = g$ , then  $h(\{a\}) = \{g\}$ . Now let  $L_1 \subseteq A_1^*$  and  $L_2 \subseteq A_2^*$  be regular languages, and let  $h_1$  and  $h_2$  be interpretations of  $A_1$  and  $A_2$ , respectively, such that  $h_1(L_1)$  and  $h_2(L_2)$  are graph languages in  $\text{Int}(\text{REG})$ . Obviously, by a renaming of symbols, we may assume that  $A_1$  and  $A_2$  are disjoint. Let  $h = h_1 \cup h_2$  be the interpretation of  $A_1 \cup A_2$  that extends both  $h_1$  and  $h_2$ . It is easy to verify that (with the appropriate conditions on types)  $h_1(L_1) \cup h_2(L_2) = h(L_1 \cup L_2)$ ,  $h_1(L_1) \circ h_2(L_2) = h(L_1 \cdot L_2)$ , and  $h_1(L_1)^+ = h(L_1^+)$ , which shows that these graph languages are also in  $\text{Int}(\text{REG})$ .

To show that  $\text{Int}(\text{REG}) \subseteq \text{REX}(\cup, \circ, +, \text{SING})$ , we first note that, since an interpretation is undefined for the empty string,  $\text{Int}(\text{REG}) = \text{Int}(\text{REG} - \lambda)$ , where  $\text{REG} - \lambda = \{L - \{\lambda\} \mid L \in \text{REG}\}$  is the class of all  $\lambda$ -free regular languages. It is well known (and easy to prove) that  $\text{REG} - \lambda$  is the smallest class of languages containing the empty language and all languages  $\{a\}$  where  $a$  is a symbol, and closed under the operations union, concatenation, and plus. By induction on this characterization we show that for every language  $L \in \text{REG} - \lambda$  and every interpretation  $h$  of the alphabet of  $L$ , if  $h(w)$  is defined for every  $w \in L$ , and  $h(L)$  is a graph language, then  $h(L) \in \text{REX}(\cup, \circ, +, \text{SING})$ . Note that by Lemma 15 (and the fact that  $\text{REG}$  is closed under intersection) we can indeed assume that  $h$  is defined for all strings in  $L$ . The inductive proof is as follows. If  $L$  is empty, then so is  $h(L)$ . If  $L = \{a\}$ , then  $h(L)$  is a singleton. If  $L = L_1 \cup L_2$ , then  $h(L) = h(L_1) \cup h(L_2)$ . Now let  $L = L_1 \cdot L_2$  and assume that  $L_1$  and  $L_2$  are nonempty (otherwise  $L$  is empty). Since, by assumption,  $h(L_1 \cdot L_2)$  is a graph language and  $h(w)$  is defined for every  $w \in L_1 \cdot L_2$ ,  $h(L_1)$  and  $h(L_2)$  are also graph languages; for  $h(L_1)$  this is proved as follows: if  $w_1, w'_1 \in L_1$ , then, for any  $w_2 \in L_2$ ,  $h(w_1 \cdot w_2) = h(w_1) \circ h(w_2)$  and similarly for  $w'_1$ , and so  $|\text{begin}(h(w_1))| = |\text{begin}(h(w_1 \cdot w_2))| = |\text{begin}(h(w'_1 \cdot w_2))| = |\text{begin}(h(w'_1))|$  and  $|\text{end}(h(w_1))| = |\text{begin}(h(w_2))| = |\text{end}(h(w'_1))|$ . Hence  $h(L) = h(L_1 \cdot L_2) = h(L_1) \circ h(L_2)$ . Finally, let  $L = L_1^+$ . Then  $h(L_1)$  is a graph language of some type  $(k, k)$  by an argument similar to the one above, and  $h(L) = h(L_1)^+$ .  $\square$

This result holds in fact for sets of morphisms of arbitrary categories (instead of the category  $\text{GR}$  of graphs, cf. Lemma 6). It generalizes a well-known characterization of the rational subsets of a monoid (see, e.g., Proposition III.2.2 of [Ber]).

The characterization of Theorem 18 still holds after adding the sum operation, extended to graph languages in the usual way: for arbitrary graph languages  $L$  and  $L'$ ,  $L \oplus L' = \{g \oplus g' \mid g \in L, g' \in L'\}$ . In other words,  $\text{Int}(\text{REG}) = \text{REX}(\cup, \circ, *, \oplus, \text{SING})$ , the smallest class of graph languages containing the empty graph language and all singleton graph languages, and closed under the operations union, concatenation, star, and sum. This is because of the following simple reason.

**Lemma 19.** *For every class of languages  $K$ , if  $\text{Int}(K)$  is closed under concatenation, then it is closed under sum.*

*Proof.* We first show that if  $M$  is in  $\text{Int}(K)$  then so is  $M \oplus \{\text{id}_k\}$  for every  $k$ . This was shown for  $\text{Val}(\text{RLIN-CFG}(\text{CS}))$  in the proof of Theorem 13, and the following argument is the same as the one used there. Let  $M = h(L)$  for some  $L \in K$  and some interpretation  $h$  of the alphabet  $A$  of  $L$ . Define  $h'(a) = h(a) \oplus \text{id}_k$  for every  $a \in A$ . Then  $h'(a_1 \cdots a_n) =$

$$(h(a_1) \oplus \text{id}_k) \circ \cdots \circ (h(a_n) \oplus \text{id}_k) = (h(a_1) \circ \cdots \circ h(a_n)) \oplus (\text{id}_k \circ \cdots \circ \text{id}_k)$$

because of strict monoidality (Lemma 6(5)), and the last expression is equal to  $h(a_1 \cdots a_n) \oplus \text{id}_k$ . This implies that  $h'(L) = h(L) \oplus \{\text{id}_k\} = M \oplus \{\text{id}_k\}$ . Similarly it can be shown that  $\{\text{id}_k\} \oplus M$  is in  $\text{Int}(K)$ .

Now, for arbitrary graph languages  $M$  and  $M'$  with  $\text{type}(M) = (m, n)$  and  $\text{type}(M') = (m', n')$ ,  $M \oplus M' = (M \circ \{\text{id}_n\}) \oplus (\{\text{id}_{m'}\} \circ M') = (M \oplus \{\text{id}_{m'}\}) \circ (\{\text{id}_n\} \oplus M')$  by strict monoidality. Hence, by the above, and the fact that  $\text{Int}(K)$  is closed under  $\circ$ ,  $M \oplus M'$  is in  $\text{Int}(K)$ .  $\square$

If we allow  $\oplus$  in our regular expressions, then, as should be clear from Theorem 7, we do not need all singleton graph languages to start with, but only those that contain elementary graphs (i.e., graphs that belong to some  $\text{EL}(\Sigma)$ , as defined in Section 3). Recall that a graph is elementary if it is an atom or one of the graphs  $I_{0,1}$ ,  $I_{1,0}$ ,  $I_{1,2}$ ,  $I_{2,1}$ ,  $\pi_{12}$ , or  $\text{id}_0$ . Let  $\text{REX}(\cup, \circ, *, \oplus, \text{ELSING})$  denote the smallest class of graph languages containing the empty graph language and all singleton graph languages with an elementary graph as element, and closed under the operations union, concatenation, star, and sum.

**Theorem 20.**  $\text{Int}(\text{REG}) = \text{REX}(\cup, \circ, *, \oplus, \text{ELSING})$ .

Note that this result is closer to the corresponding result for regular languages, for which only singleton languages  $\{a\}$  are needed where  $a$  is a symbol.

The next class  $K$  of interest is the class CF of context-free languages. We will show that  $\text{Int}(\text{CF})$  is a (proper) subclass of HR, the class of graph languages generated by HR grammars. In fact, it is rather obvious that  $\text{Int}(\text{CF})$  is exactly the class of languages generated by context-free graph grammars over CS that do not use the sum operation. Thus,  $\text{Int}(\text{CF})$  is the class of equational subsets of the algebra of graphs with the concatenation operation. Inclusion in HR then follows from Theorem 11. As in the previous theorems, we have to cope with the technical trouble of typing.

Let  $\text{CS}_\circ = \text{CS} - \{\oplus\} = \{\circ\} \cup \{c_g \mid g \in \text{GR}\}$ . If  $\oplus$  does not occur in the productions of a context-free graph grammar over CS, then we also say that it is over  $\text{CS}_\circ$ . By  $\text{Val}(\text{CFG}(\text{CS}_\circ))$  we denote the class of all graph languages generated by context-free graph grammars over  $\text{CS}_\circ$ . Note that, by definition,  $\text{Val}(\text{RLIN-CFG}(\text{CS})) \subseteq \text{Val}(\text{CFG}(\text{CS}_\circ))$ .

**Theorem 21.**  $\text{Int}(\text{CF}) = \text{Val}(\text{CFG}(\text{CS}_\circ))$ .

*Proof.* We first show that  $\text{Val}(\text{CFG}(\text{CS}_\circ)) \subseteq \text{Int}(\text{CF})$ . Let  $G = (N, T, P, S)$  be a context-free graph grammar over  $\text{CS}_\circ$ . Every production of  $G$  is of the form  $X \rightarrow \alpha_1 \circ \alpha_2 \circ \dots \circ \alpha_k$  with  $k \geq 1$  and  $\alpha_i \in N$  or  $\alpha_i = c_g$  for some  $g \in \text{GR}$ . Note that we can drop the parentheses from the right-hand sides, due to the associativity of concatenation. Thus,  $G$  generates expressions of the form  $c_{g_1} \circ \dots \circ c_{g_n}$  with  $n \geq 1$ , and  $\text{val}(L(G)) = \{g_1 \circ \dots \circ g_n \mid c_{g_1} \circ \dots \circ c_{g_n} \in L(G)\}$ .

Define the (ordinary) context-free grammar  $G' = (N, T', P', S)$  where  $T'$  is the set of all  $c_g$  in  $T$ , and  $P' = \{X \rightarrow \alpha_1 \alpha_2 \dots \alpha_k \mid X \rightarrow \alpha_1 \circ \alpha_2 \circ \dots \circ \alpha_k \in P\}$ . Obviously,  $L(G') = \{c_{g_1} \dots c_{g_n} \mid c_{g_1} \circ \dots \circ c_{g_n} \in L(G)\}$ . Now let  $h : T' \rightarrow \text{GR}$  be the interpretation such that  $h(c_g) = g$ . Then, clearly,  $\text{val}(L(G)) = h(L(G'))$ . Hence  $\text{val}(L(G))$  is in  $\text{Int}(\text{CF})$ .

We now show that  $\text{Int}(\text{CF}) \subseteq \text{Val}(\text{CFG}(\text{CS}_\circ))$ . By Lemma 15 (and the fact that  $\text{CF}$  is closed under intersection with regular languages), every  $\text{Int}(\text{CF})$  graph language is of the form  $h(L)$  where  $L$  is a context-free language such that  $h$  is defined for every string in  $L$ . In particular,  $L$  is  $\lambda$ -free. Let  $G = (N, T, P, S)$  be a context-free grammar generating  $L$ . We may assume that the right-hand sides of the productions of  $G$  are non-empty. Define the context-free grammar  $G' = (N, T', P', S)$  such that  $T' = \{\circ\} \cup \{c_{h(a)} \mid a \in T\}$ , and  $P' = \{X \rightarrow \psi(\alpha_1) \circ \psi(\alpha_2) \circ \dots \circ \psi(\alpha_k) \mid X \rightarrow \alpha_1 \alpha_2 \dots \alpha_k \in P\}$ , where  $\psi(a) = c_{h(a)}$  for every  $a \in T$ , and  $\psi(X) = X$  for every  $X \in N$ . Obviously,  $L(G') = \{c_{h(a_1)} \circ \dots \circ c_{h(a_n)} \mid a_1 \dots a_n \in L(G)\}$ , and so  $\text{val}(L(G')) = h(L(G)) = h(L)$ .

The only thing that remains to be proved (and this is the “technical trouble”) is that  $G'$  is a context-free graph grammar over  $\text{CS}$  (and it obviously is over  $\text{CS}_\circ$ ). In other words, we have to turn  $N$  into a typed alphabet, such that the right-hand sides of the productions are expressions with the same type as the left-hand sides. To this aim we investigate the grammar  $G$  in more detail.

We claim that for all strings  $w \in T^*$  generated by a given nonterminal  $X$  of  $G$ ,  $h(w)$  is defined and  $\text{type}(h(w))$  is the same for all such  $w$ . Here we will use the fact that  $h$  is defined for all strings in  $L$ . The proof is similar to the argument used at the end of the proof of Theorem 18. Let  $X \Rightarrow^* w_1$  and  $X \Rightarrow^* w_2$ , with  $w_i \in T^*$ . Consider some  $u, v \in T^*$  such that  $S \Rightarrow^* uXv$  (assuming  $G$  to be reduced). Then  $uw_1v, uw_2v \in L$ . We now show that  $|\text{begin}(h(w_1))| = |\text{begin}(h(w_2))|$ . Let  $\text{type}(L) = (m, n)$ . If  $u = \lambda$ , then  $w_i v \in L$  and  $|\text{begin}(h(w_i))| = |\text{begin}(h(w_i v))| = m$ . Note that if  $h$  is defined for a string  $w$ , then it is also defined for every non-empty substring of  $w$ . Now let  $u \neq \lambda$ . Since  $h$  is defined on  $uw_i v$ ,  $|\text{begin}(h(w_i))| = |\text{end}(h(u))|$ . In the same way it can be shown that  $|\text{end}(h(w_1))| = |\text{end}(h(w_2))|$ .

We turn  $N$  into a typed alphabet by defining  $\text{type}(X) = \text{type}(h(w))$  if  $X \Rightarrow^* w$  in  $G$ . We now have to show that for every production  $X \rightarrow \alpha_1 \dots \alpha_k$  of  $G$ ,  $\psi(\alpha_1) \circ \dots \circ \psi(\alpha_k)$  is a well-formed expression over  $\text{CS}$  and  $N$ , of the same type as  $X$ . Note first that for every  $\alpha \in N \cup T$  and every  $w \in T^*$ , if  $\alpha \Rightarrow^* w$ , then  $h(w)$  is defined and  $\text{type}(\psi(\alpha)) = \text{type}(h(w))$ . Now consider  $\alpha_1, \dots, \alpha_k$  in  $N \cup T$ , and let  $\alpha_i \Rightarrow^* w_i \in T^*$ , for  $1 \leq i \leq k$ . Then  $\psi(\alpha_1) \circ \dots \circ \psi(\alpha_k)$  is a well-formed expression (over  $\text{CS}$  and  $N$ ) of type  $(m, n)$  if and only if  $h(w_1 \dots w_k)$  is defined and  $\text{type}(h(w_1 \dots w_k)) = (m, n)$ . Consider the derivation  $S \Rightarrow^* uXv \Rightarrow u\alpha_1 \dots \alpha_k v \Rightarrow^* uw_1 \dots w_k v = z$ . Since  $z \in L$ ,  $h(z)$  is defined, and so  $h(w_1 \dots w_k)$

is defined. Moreover, since  $X \Rightarrow^* w_1 \cdots w_k$ ,  $\text{type}(h(w_1 \cdots w_k)) = \text{type}(X)$ . This shows that  $\psi(\alpha_1) \circ \cdots \circ \psi(\alpha_k)$  is a well-formed expression over CS and  $N$ , of the same type as  $X$ , as required.  $\square$

**Theorem 22.**  $\text{Int}(\text{CF}) \subset \text{HR}$ .

*Proof.* Inclusion follows immediately from Theorems 21 and 11. Proper inclusion is a consequence of Corollary 17: the set of all trees is in HR, but is not of bounded pathwidth, as can easily be seen (for a characterization of the trees of pathwidth  $k$ , see [ELLST]).  $\square$

Since REG is closed under intersection, the proof of Theorem 21 also works for  $\text{Int}(\text{REG})$ . In fact, the proof preserves the right-linearity of the grammars. Hence,  $\text{Int}(\text{REG}) = \text{Val}(\text{RLIN-CFG}(\text{CS}))$ . As an example, the graph language of clothes lines is generated by the right-linear context-free graph grammar with productions  $S \rightarrow c_{h(a)} \circ X$ ,  $X \rightarrow c_{h(b)} \circ X$ ,  $X \rightarrow c_{h(c)} \circ X$ , and  $X \rightarrow c_{h(d)}$ , where  $h(a), h(b), h(c), h(d)$  are the graphs in Fig. 4,  $\text{type}(S) = (0, 0)$ , and  $\text{type}(X) = (1, 0)$ .

Together with Theorem 13, this shows that the graph languages that are interpretations of a regular language are precisely those that can be generated by linear HR grammars.

**Theorem 23.**  $\text{Int}(\text{REG}) = \text{LIN-HR}$ .

Similarly, the proof of Theorem 21 preserves linearity of the grammars (and LIN is closed under intersection with regular languages). Since  $\text{Val}(\text{LIN-CFG}(\text{CS}_\circ))$  is inbetween  $\text{Val}(\text{RLIN-CFG}(\text{CS}))$  and  $\text{Val}(\text{LIN-CFG}(\text{CS}))$ , we obtain the next result by Corollary 12 and Theorem 13.

**Theorem 24.**  $\text{Int}(\text{LIN}) = \text{LIN-HR}$ .

The results of this section suggest that for graph languages, regularity and linearity are the same. We have a class of graph languages that may be called the class of regular graph languages on the one hand (because it is equal to  $\text{Int}(\text{REG})$  and to  $\text{REX}(\cup, \circ, *, \text{SING})$ ), and may be called the class of linear graph languages on the other hand (because it is equal to  $\text{Int}(\text{LIN})$  and to  $\text{LIN-HR}$ ).

It will be shown in Section 7 that  $\text{Int}(\text{REG})$  is a proper subclass of  $\text{Int}(\text{CF})$ .

## 6 Characterizations of $\text{Int}(K)$

In this section and the next, we investigate properties of the class of graph languages  $\text{Int}(K)$  for arbitrary classes of string languages  $K$ . However, to avoid trivialities we will mainly be interested in classes  $K$  that are *closed under sequential machine mappings*, where a sequential machine is a transducer which works like an ordinary nondeterministic finite automaton that, moreover, at each step outputs one symbol (thus it is a special case of the generalized sequential machine, or *gsm*, which outputs a string at each step, see, e.g., [HopUll]).

Equivalently,  $K$  is closed under intersection with regular languages and under *alphabetical substitutions* (where an alphabetical substitution from alphabet  $A$  to alphabet  $B$  is a relation  $\rho \subseteq A \times B$  that is extended to a function from  $A^*$  to the finite subsets of  $B^*$  by  $\rho(a_1 \cdots a_n) = \{b_1 \cdots b_n \mid (a_i, b_i) \in \rho, 1 \leq i \leq n\}$ ).

In this section we present two characterizations of  $\text{Int}(K)$ . The first characterization of  $\text{Int}(K)$  is through the notion of replacement, as defined in Section 2.2. The following closure property of  $\text{Int}(K)$  will be useful.

**Lemma 25.** *For every class  $K$ ,  $\text{Int}(K)$  is closed under replacements.*

*Proof.* Clearly,  $\phi(h(L)) = h'(L)$ , where  $h'$  is defined by  $h'(a) = \phi(h(a))$  for every  $a \in A$ . In fact, for  $a_1, \dots, a_n \in A$ ,  $\phi(h(a_1 \cdots a_n)) = \phi(h(a_1) \circ \cdots \circ h(a_n)) = \phi(h(a_1)) \circ \cdots \circ \phi(h(a_n)) = h'(a_1) \circ \cdots \circ h'(a_n) = h'(a_1 \cdots a_n)$ , where we have used Lemma 4(1).  $\square$

The characterization of  $\text{Int}(K)$  is based on the fact that every interpretation can be decomposed into an “atomic” interpretation and a replacement. An interpretation  $h : A \rightarrow \text{GR}$  of  $A$  is *atomic* if  $A$  is a typed alphabet and  $h(a) = \text{atom}(a)$  for every  $a \in A$ . By  $\text{AtInt}(K)$  we denote the set of all  $h(L) \in \text{Int}(K)$  such that  $h$  is an atomic interpretation. Recall that  $\text{Repl}$  denotes the class of all replacements.

**Theorem 26.** *For every class  $K$ ,  $\text{Int}(K) = \text{Repl}(\text{AtInt}(K))$ .*

*Proof.* One inclusion follows from Lemma 25. For the other inclusion, let  $h : A \rightarrow \text{GR}$  be an interpretation. Turn  $A$  into a typed alphabet by defining  $\text{type}(a) = \text{type}(h(a))$  for every  $a \in A$ . Let  $t$  be the unique atomic interpretation of the typed alphabet  $A$ , and let  $\phi : A \rightarrow \text{GR}$  be the replacement defined by  $\phi(a) = h(a)$  for every  $a \in A$  (i.e.,  $\phi$  is  $h$  viewed as a replacement). Clearly, for every string  $w \in A^*$ ,  $\phi(t(w)) = h(w)$ . In fact, if  $w = a_1 \cdots a_n$ , then  $\phi(t(w)) = \phi(t(a_1) \circ \cdots \circ t(a_n)) = \phi(t(a_1)) \circ \cdots \circ \phi(t(a_n)) = \phi(a_1) \circ \cdots \circ \phi(a_n) = h(a_1) \circ \cdots \circ h(a_n) = h(w)$ , by Lemma 4(1) and because  $\phi(\text{atom}(a)) = \phi(a)$  for every  $a \in A$ .  $\square$

Note that Lemma 25 and Theorem 26 together show that  $\text{Int}(K)$  is the smallest class of graph languages containing  $\text{AtInt}(K)$  and closed under replacements.

The second characterization of  $\text{Int}(K)$  generalizes the characterization of  $\text{Int}(\text{REG})$  in Theorem 23. To this aim we consider *controlled* linear HR grammars, in the obvious sense. Let  $G = (N, T, P, S)$  be an HR grammar, and let  $C$  be a string language over  $P$  (where  $P$  is viewed as an alphabet). The graph language generated by  $G$  under control  $C$  is the set of all graphs  $g \in \text{GR}(T)$  for which there is a derivation  $\text{atom}(S) \Rightarrow_{p_1} g_1 \Rightarrow_{p_2} g_2 \cdots \Rightarrow_{p_n} g_n$  with  $g_n = g$ , such that the string  $p_1 p_2 \cdots p_n$  is in  $C$ . Recall that  $\Rightarrow_p$  denotes a derivation step of  $G$  that uses production  $p$ . Thus, the control language  $C$  specifies the sequences of productions that the grammar  $G$  is allowed to use in its derivations. If  $C$  is taken from a class of languages  $K$ , the grammar  $G$  together with the control language  $C$  is also called a  *$K$ -controlled HR grammar*.

For a class  $K$  of string languages, we denote by  $\text{LIN-HR}(K)$  the class of graph languages generated by  $K$ -controlled linear HR grammars. Generalizing Theorem 23 and its proof we obtain the next result.



**Theorem 27.** *For every class  $K$  that is closed under sequential machine mappings,  $\text{Int}(K) = \text{LIN-HR}(K)$ .*

*Proof.* In what follows we assume, without loss of generality, that all languages in  $K$  are  $\lambda$ -free (if  $K' = \{L - \{\lambda\} \mid L \in K\}$ , then  $\text{Int}(K') = \text{Int}(K)$ ,  $\text{LIN-HR}(K') = \text{LIN-HR}(K)$ ,  $K' = \{L \in K \mid \lambda \notin L\}$  because a sequential machine mapping can be used to remove  $\lambda$ , and  $K'$  is closed under sequential machine mappings, because  $K$  is closed under sequential machine mappings and sequential machine mappings are length-preserving).

With analogous definitions as above, we can define the controlled versions of context-free grammars and of context-free graph grammars over CS. We first show that  $K$  is the class of languages generated by the  $K$ -controlled right-linear context-free grammars. In one direction, let  $L \in K$  with  $L \subseteq T^*$ . Define the right-linear context-free grammar  $G = (\{S\}, T, P, S)$ , where  $P$  consists of all productions  $p_a : S \rightarrow aS$  and  $p'_a : S \rightarrow a$  for all  $a \in T$ . Let  $L'$  be the control language that is obtained from  $L$  by a sequential machine mapping that changes each string  $a_1 \cdots a_{n-1}a_n$  into the string  $p_{a_1} \cdots p_{a_{n-1}}p'_{a_n}$ . Then  $L$  is the language generated by  $G$  under control  $L'$ . In the other direction, let  $G = (N, T, P, S)$  be a right-linear context-free grammar, let  $C \in K$  be a control language, and let  $L$  be the language generated by  $G$  under control  $C$ . Let  $\phi$  be the sequential machine mapping that, for a given string  $p_1 \cdots p_n \in P^*$ , checks whether there is a derivation  $S \Rightarrow_{p_1} w_1 \cdots \Rightarrow_{p_n} w_n$  with  $w_n \in T^*$  (by simulating  $G$  in its finite control) and changes each production into the unique terminal symbol that occurs in its right-hand side. Then, clearly,  $\phi(C) = L$  and hence  $L$  is in  $K$ .

It is now easy to generalize the proof of Theorem 21 to the case of  $K$ -controlled right-linear grammars. This shows that  $\text{Int}(K)$  is equal to the  $K$ -controlled version of  $\text{Val}(\text{RLIN-CFG}(\text{CS}))$ , i.e., to the class of graph languages generated by  $K$ -controlled right-linear context-free graph grammars over CS. Thus, it now suffices to check that the proof of Theorem 13 can be generalized to the  $K$ -controlled case.

First we check that the proof of Theorem 11 can be generalized to the  $K$ -controlled case, for linear grammars. We have proved a relationship between the derivations  $X \Rightarrow^* e$  of a context-free graph grammar  $G$  and the derivations  $\text{atom}(X) \Rightarrow^* \text{val}(e)$  of an HR grammar  $G'$ . This proof can be extended to show that if the sequence of productions used in the first derivation is  $p_1 p_2 \cdots p_n$ , then the sequence of productions used in the second derivation is  $p'_1 p'_2 \cdots p'_n$ , where, for each production  $p = X \rightarrow t$  of  $G$ ,  $p'$  is the production  $X \rightarrow \text{val}(t)$  of  $G'$ ; to see this, note that in the linear case of Proposition 2, the sequence of productions used in  $h \Rightarrow^* g$  equals the sequence of productions used in  $\text{atom}(Y) \Rightarrow^* \phi(Y)$ , where  $Y$  is the unique nonterminal in  $\text{lab}(h)$ . Thus, if  $C$  is the control language of  $G$ , then the control language for  $G'$  can be obtained from  $C$  by a sequential machine mapping that changes every  $p$  into  $p'$ .

We now consider the proof of Theorem 13. Clearly, for the LIN-HR grammar  $\overline{G}$  we can take the same control language that is used by  $G$  (modulo a renaming), because there is a clear one-to-one correspondence between their productions. In the remaining two constructions, we can take the same control language in the

first case, and in the second case we can apply a sequential machine mapping to the control language that changes the first production of a production sequence into the corresponding production for the new initial nonterminal  $S'$  (which was the reason to use that particular construction).  $\square$

## 7 Comparison of $\text{Int}(K)$ and $\text{Int}(K')$

From Theorems 23 and 24 we know that  $\text{Int}(\text{REG}) = \text{Int}(\text{LIN}) = \text{LIN-HR}$ . It can be shown by a direct construction (see [Ver]) that even  $\text{Int}(\text{DB}) = \text{LIN-HR}$ , where DB is the class of derivation bounded context-free languages (see, e.g., Section VI.10 of [Sal], where they are called languages of “finite index”). One now wonders when  $\text{Int}(K) = \text{Int}(K')$ , and in particular one wonders how much larger the class  $K$  can be made without enlarging the class  $\text{Int}(K)$ .

It is easy to see that for every given class  $K$  there is a largest class  $K'$  such that  $\text{Int}(K') = \text{Int}(K)$ . We will call this the *extension of  $K$* , denoted  $\text{Ext}(K)$ . In fact,  $\text{Ext}(K) = \bigcup \{K' \mid \text{Int}(K') = \text{Int}(K)\}$ . Note that, for arbitrary classes  $K$  and  $K'$ ,  $\text{Int}(K) = \text{Int}(K')$  if and only if  $\text{Ext}(K) = \text{Ext}(K')$ .

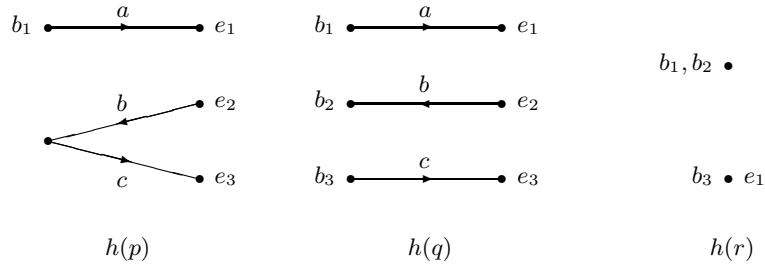


Fig. 6. Interpretation  $h$ .

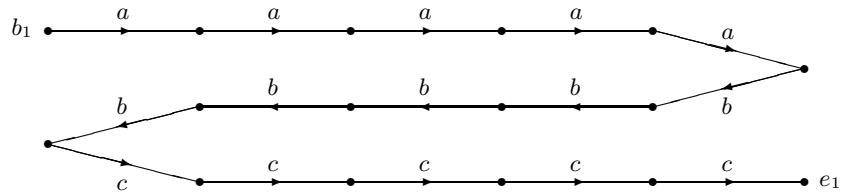


Fig. 7. The graph  $\text{gr}(a^5 b^5 c^5) = h(pq^4 r)$ .

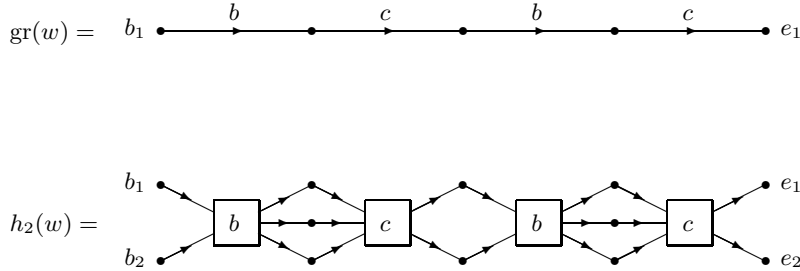
In the next theorem we give a characterization of  $\text{Ext}(K)$ . For a class  $G$  of graph languages, let  $\text{Str}(G)$  denote the class of string languages  $L$  such that  $\text{gr}(L)$  is in  $G$ . Here,  $\text{gr}(L) = \{\text{gr}(w) \mid w \in L\}$ , and, for a string  $w = a_1 \cdots a_n$  with  $n \geq 1$ ,  $\text{gr}(w) = (V, E, s, t, l, \text{begin}, \text{end})$  is the (ordinary) graph of type  $(1, 1)$  with  $V = \{x_1, \dots, x_{n+1}\}$ ,  $E = \{e_1, \dots, e_n\}$ ,  $s(e_i) = x_i$ ,  $t(e_i) = x_{i+1}$ , and  $l(e_i) = a_i$  for every  $1 \leq i \leq n$ ,  $\text{begin} = x_1$  and  $\text{end} = x_{n+1}$ . Thus,  $\text{gr}(w)$  encodes  $w$  in the obvious way: it is a path with the symbols of  $w$  as edge labels.

As a classical example, the language  $L = \{a^n b^n c^n \mid n \geq 1\}$  is in  $\text{Str}(\text{Int}(\text{REG}))$ , because  $L = h(M)$  with  $M = pq^*r$  and  $h$  is shown in Fig. 6. The graph  $\text{gr}(a^5 b^5 c^5) = h(pq^4 r)$  is shown in Fig. 7.

Note that  $\text{gr}(\lambda)$  is not defined. This implies that  $L \in \text{Str}(G)$  iff  $L - \{\lambda\} \in \text{Str}(G)$ . It also implies that ‘gr’ is the unique atomic interpretation which is obtained by viewing every symbol as having type  $(1, 1)$ ; hence  $\text{gr}(L) \in \text{Int}(K)$  for every  $L \in K$ , which means that  $K \subseteq \text{Str}(\text{Int}(K))$ .

**Theorem 28.** *For every class  $K$  that is closed under sequential machine mappings,  $\text{Ext}(K) = \text{Str}(\text{Int}(K))$ .*

*Proof.* Clearly, if  $\text{Int}(K') = \text{Int}(K)$ , then  $K' \subseteq \text{Str}(\text{Int}(K')) = \text{Str}(\text{Int}(K))$ . Thus, it remains to show that  $\text{Int}(\text{Str}(\text{Int}(K))) = \text{Int}(K)$ . Since  $K \subseteq \text{Str}(\text{Int}(K))$ ,  $\text{Int}(K) \subseteq \text{Int}(\text{Str}(\text{Int}(K)))$ . For the other inclusion it suffices, by Lemma 25 and Theorem 26, to show that  $\text{AtInt}(\text{Str}(\text{Int}(K))) \subseteq \text{Int}(K)$ . To prove this, let  $L_1 \in K$ , let  $h_1$  be an interpretation of the alphabet  $A$  of  $L_1$  such that  $h_1(L_1) = \text{gr}(L_2)$  for some  $\lambda$ -free string language  $L_2$ , and let  $h_2$  be an atomic interpretation of the alphabet  $B$  of  $L_2$ . Thus,  $h_1 : A \rightarrow \text{GR}(B)$  where each symbol from  $B$  has type  $(1, 1)$ . However, for the atomic interpretation  $h_2$  each symbol  $b$  from  $B$  has another (arbitrary) type that we will denote by  $\text{type}(b)$ . Note that  $h_2(b) = \text{atom}(b)$ , where  $\text{type}(\text{atom}(b)) = \text{type}(b)$ ; hence  $\text{type}(b) = (|\text{begin}(h_2(b))|, |\text{end}(h_2(b))|)$ .



**Fig. 8.** Graphs  $\text{gr}(w)$  and  $h_2(w)$  for  $w = bcbc$ , with  $\text{type}(b) = (2, 3)$  and  $\text{type}(c) = (3, 2)$ .

We have to construct a language  $L \in K$  and an interpretation  $h$  such that  $h(L) = h_2(L_2)$ . For a string  $w \in L_2$  such that  $h_2(w)$  is defined, the graph  $h_2(w)$

can be obtained from the graph  $\text{gr}(w)$  (which is an element of  $h_1(L_1)$ ) in an easy way, as follows (see Fig. 8). Each node  $v$  of  $\text{gr}(w)$  has to be “expanded” into a sequence of distinct nodes  $(v, 1), (v, 2), \dots, (v, \mu(v))$ , where  $\mu$  stands for “multiplicity”. Clearly,  $\mu(v)$  is determined by  $\text{type}(b)$ , where  $b$  is the label of an edge  $e$  incident with  $v$ : if  $e$  enters  $v$ , then  $\mu(v) = |\text{end}(h_2(b))|$ , and if  $e$  leaves  $v$ , then  $\mu(v) = |\text{begin}(h_2(b))|$ . Every edge  $e$  of  $\text{gr}(w)$ , with source  $u$  and target  $v$ , should be replaced by an edge  $e$  with sources  $(u, 1), \dots, (u, \mu(u))$  and targets  $(v, 1), \dots, (v, \mu(v))$  (and the same label). In Fig. 8, the multiplicity of the nodes of  $\text{gr}(w)$  is 2, 3, 2, 3, 2, respectively. The edges of  $h_2(w)$  are drawn as squares, with “tentacles” from their sources and to their targets (where we assume that the tentacles are ordered, e.g., from top to bottom).

We now define this expansion process formally, for arbitrary graphs in  $\text{GR}(B)$ . Let  $M$  be a number such that for all  $b \in B$ , if  $\text{type}(b) = (m, n)$ , then  $m, n \leq M$ . A *decoration* of a graph  $g \in \text{GR}(B)$  is a mapping  $\mu : V_g \rightarrow \{1, \dots, M\}$  such that for every edge  $e$  of  $g$  with  $s(e) = u$ ,  $t(e) = v$ , and  $l(e) = b$ :  $\text{type}(b) = (\mu(u), \mu(v))$ , i.e.,  $\mu(u) = |\text{begin}(h_2(b))|$  and  $\mu(v) = |\text{end}(h_2(b))|$ . Note that for every graph  $\text{gr}(w)$ ,  $h_2(w)$  is defined if and only if there is a decoration  $\mu$  of  $\text{gr}(w)$ , and in that case  $\mu$  is unique. For a graph  $g \in \text{GR}(B)$  and a decoration  $\mu$  of  $g$ , the *expansion* of  $g$  by  $\mu$ , denoted  $\text{exp}(g, \mu)$ , is the graph that has all nodes  $(v, i)$  where  $v$  is a node of  $g$  and  $1 \leq i \leq \mu(v)$ ; it has the same edges as  $g$  (with the same labels), but if  $s(e) = u$  and  $t(e) = v$  in  $g$ , then  $s(e) = (u, 1) \cdots (u, \mu(u))$  and  $t(e) = (v, 1) \cdots (v, \mu(v))$  in  $\text{exp}(g, \mu)$ ; finally, if  $\text{begin}(g) = v_1 v_2 \cdots v_k$ , then  $\text{begin}(\text{exp}(g, \mu)) =$

$$(v_1, 1) \cdots (v_1, \mu(v_1)) \cdots (v_2, 1) \cdots (v_2, \mu(v_2)) \cdots (v_k, 1) \cdots (v_k, \mu(v_k)),$$

and similarly for  $\text{end}(g)$  and  $\text{end}(\text{exp}(g, \mu))$ . It should be clear that for every  $w \in L_2$  for which  $h_2(w)$  is defined,  $h_2(w) = \text{exp}(\text{gr}(w), \mu)$  where  $\mu$  is the unique decoration mentioned above.

Based on this idea of expansion we now change  $L_1$  into  $L$  and  $h_1$  into  $h$ , as follows. Let  $A'$  be the alphabet consisting of all pairs  $(a, \mu)$  with  $a \in A$  and  $\mu$  is a decoration of  $h_1(a)$ . Intuitively,  $\mu$  is a guess of the multiplicities of the nodes of  $h_1(a)$  as they will occur in a graph of  $h_1(L_1)$ ; for nodes that are incident with an edge, this multiplicity is determined by the label of that edge, but for the other nodes (which are necessarily begin or end nodes because  $\text{gr}(w)$  has no isolated nodes) their multiplicity will only be clear after concatenation. Let  $\rho$  be the alphabetical substitution that substitutes all possible  $(a, \mu)$  for  $a$ , i.e.,  $\rho = \{(a, (a, \mu)) \mid (a, \mu) \in A'\}$ . Thus,  $\rho(L_1) = \{(a_1, \mu_1) \cdots (a_n, \mu_n) \mid a_1 \cdots a_n \in L_1, (a_j, \mu_j) \in A'\}$ . Let  $R$  be the regular language over  $A'$  that consists of all strings  $(a_1, \mu_1) \cdots (a_n, \mu_n)$  such that  $h_1(a_1 \cdots a_n)$  is defined and  $\mu_j(\text{end}(h_1(a_j))(i)) = \mu_{j+1}(\text{begin}(h_1(a_{j+1}))(i))$  for all relevant  $j$  and  $i$  (to be precise: for all  $1 \leq j < n$  and all  $1 \leq i \leq |\text{end}(h_1(a_j))|$ ; note that  $|\text{end}(h_1(a_j))| = |\text{begin}(h_1(a_{j+1}))|$  because  $h_1(a_1 \cdots a_n)$  is defined). In words, the language  $R$  checks that the guessed multiplicity of the  $i$ th end node of  $h_1(a_j)$  equals that of the  $i$ th begin node of  $h_1(a_{j+1})$ . Thus,  $R$  checks that the guessed multiplicities are consistent with the identification of nodes when concatenating the graphs

$h_1(a_1), \dots, h_1(a_n)$  (and hence are the correct multiplicities). It should be clear that  $R$  is indeed regular (cf. Lemma 15). We now define  $L = \rho(L_1) \cap R$ ; since, by assumption,  $K$  is closed under sequential machine mappings,  $L$  is in  $K$ . Finally, for  $(a, \mu) \in A'$  we define  $h(a, \mu) = \exp(h_1(a), \mu)$ .

It should now be clear that  $h(L) = h_2(L_2)$ . A formal proof can be given as follows.

Let a *decorated graph* be a pair  $(g, \mu)$  with  $g \in \text{GR}(B)$  and  $\mu$  a decoration of  $g$ . Define the concatenation of decorated graphs, as follows. For decorated graphs  $(g_1, \mu_1)$  and  $(g_2, \mu_2)$ , if  $g_1 \circ g_2$  is defined and  $\mu_1(\text{end}(g_1)(i)) = \mu_2(\text{begin}(g_2)(i))$  for all  $i$ , then their concatenation is  $(g_1, \mu_1) \circ (g_2, \mu_2) = (g_1 \circ g_2, \mu)$  with  $\mu([x]_R) = \mu_j(x)$  if  $x \in V_{g_j}$  (where we assume the terminology of Definition 3). It is easy to see that, the other way around, if  $g_1 \circ g_2$  is defined and  $(g_1 \circ g_2, \mu)$  is a decorated graph, then there exist decorations  $\mu_1$  and  $\mu_2$  of  $g_1$  and  $g_2$ , respectively, such that  $(g_1, \mu_1) \circ (g_2, \mu_2) = (g_1 \circ g_2, \mu)$ . As a basic property of ‘exp’ it can be shown that it is a homomorphism with respect to the concatenation of decorated graphs:  $\exp((g_1, \mu_1) \circ (g_2, \mu_2)) = \exp(g_1, \mu_1) \circ \exp(g_2, \mu_2)$ .

Now consider some string  $a_1 \cdots a_n \in L_1$ . Then  $h_2(\text{gr}^{-1}(h_1(a_1 \cdots a_n)))$  is defined if and only if there exist decorations  $\mu_i$  such that  $(h_1(a_i), \mu_i)$  is a decorated graph, for all  $i$ , and their concatenation  $(h_1(a_1), \mu_1) \circ \cdots \circ (h_1(a_n), \mu_n)$  is defined. And this is if and only if there exist  $\mu_i$  such that  $(a_i, \mu_i) \in A'$  and  $(a_1, \mu_1) \cdots (a_n, \mu_n) \in R$ . Moreover, in that case, for the unique decoration  $\mu$  of  $h_1(a_1 \cdots a_n)$ ,

$$\begin{aligned} h_2(\text{gr}^{-1}(h_1(a_1 \cdots a_n))) &= \exp(h_1(a_1 \cdots a_n), \mu) = \\ \exp((h_1(a_1), \mu_1) \circ \cdots \circ (h_1(a_n), \mu_n)) &= \\ \exp(h_1(a_1), \mu_1) \circ \cdots \circ \exp(h_1(a_n), \mu_n) &= \\ h(a_1, \mu_1) \circ \cdots \circ h(a_n, \mu_n) &= h((a_1, \mu_1) \cdots (a_n, \mu_n)). \end{aligned}$$

This shows that  $h_2(L_2) = h(L)$ . □

As a corollary of Theorem 28 we obtain that for arbitrary  $K$  and  $K'$  (both closed under sequential machine mappings),  $\text{Int}(K) = \text{Int}(K')$  if and only if  $\text{Str}(\text{Int}(K)) = \text{Str}(\text{Int}(K'))$ . This means that the graph generating power of  $K$  is completely determined by its string generating power (with strings coded as graphs by the mapping  $\text{gr}$ ).

We now show that  $\text{Ext}(K)$  is a class of languages that is well known in formal language theory. By  $2\text{DGSM}(K)$  we denote the class of images of languages from  $K$  under *2dgsms mappings*, i.e., the class of all  $f(L)$  where  $f$  is a 2dgsms mapping and  $L \in K$ . A *2dgsms* (i.e., a two-way deterministic generalized sequential machine) is a deterministic finite automaton that can move in two directions on its input tape (with endmarkers), and outputs a (possibly empty) string at each step. As an example,  $\{a^n b^n c^n \mid n \geq 1\}$  is in  $2\text{DGSM}(\text{REG})$ , because it is easy to construct a 2dgsms that translates  $pq^n r$  into  $a^{n+1} b^{n+1} c^{n+1}$  for every  $n \in \mathbf{N}$ . The proof of the next result is obtained by generalizing the proof in [EngHey] that  $\text{Str}(\text{LIN-HR})$  equals the class of output languages of 2dgsms mappings. We say

that  $K$  is *nontrivial* if it is not a subset of  $\{\emptyset, \{\lambda\}\}$ ; in other words,  $K$  contains at least one language that contains a nonempty string.

**Lemma 29.** *For every nontrivial class  $K$  that is closed under sequential machine mappings,  $\text{Str}(\text{LIN-HR}(K)) = 2\text{DGSM}(K)$ .*

*Proof.* To reduce the proof to a generalization of the proof in [EngHey], we have to deal with some technical details. In particular, the coding ‘gr’ (and hence the operation ‘Str’) is defined in a different way (see Definition 2.2 of [EngHey]). We will discuss this in steps. Note that, by Theorem 27,  $\text{Str}(\text{LIN-HR}(K)) = \text{Str}(\text{Int}(K))$ .

First of all, let  $\text{gr}_1$  be defined in the same way as gr, except that additionally  $\text{gr}_1(\lambda) = \text{id}_1$ ; and let  $\text{Str}_1$  be defined in the same way as Str, with  $\text{gr}_1$  instead of gr. We claim that  $\text{Str}_1(\text{Int}(K)) = \text{Str}(\text{Int}(K))$ . Thus, we have to show that for every language  $L$ ,  $\text{gr}_1(L) \in \text{Int}(K)$  iff  $\text{gr}(L) \in \text{Int}(K)$ . This is obvious if  $\lambda \notin L$ . If  $\lambda \in L$ , then  $\text{gr}_1(L) = \text{gr}(L) \cup \{\text{id}_1\}$ . Assume first that  $\text{gr}(L)$  is in  $\text{Int}(K)$ . If  $\text{gr}(L) = \emptyset$ , we have to show that  $\{\text{id}_1\} \in \text{Int}(K)$ . Since  $K$  is not a subset of  $\{\emptyset, \{\lambda\}\}$ ,  $K$  contains a language  $M \subseteq A^*$  such that  $M$  contains at least one nonempty string. Define the interpretation  $h$  with  $h(a) = \text{id}_1$  for all  $a \in A$ . Then  $h(M) = \{\text{id}_1\}$  (because  $\text{id}_1 \circ \text{id}_1 = \text{id}_1$ ). Now let  $\text{gr}(L) \neq \emptyset$ . Let  $\text{gr}(L) = h(M)$  for some interpretation  $h$  and some  $M \in K$ . Then  $M$  must contain a nonempty string  $w$ . Let  $b$  be a new symbol, not in the alphabet  $A$  of  $M$ , and define the interpretation  $h'$  of  $A \cup \{b\}$  such that  $h'(a) = h(a)$  for every  $a \in A$  and  $h'(b) = \text{id}_1$ . Then  $h'(M \cup \{b^{|w|}\}) = h(M) \cup \{\text{id}_1\} = \text{gr}_1(L)$ . It is easy to see that there is a sequential machine mapping that transforms  $M$  into  $M \cup \{b^{|w|}\}$ . This proves one direction of the equivalence. To show the other direction, assume that  $\text{gr}_1(L) \in \text{Int}(K)$ . If  $\text{id}_1 \in \text{gr}(L)$ , then there is nothing to prove. Now let  $\text{id}_1 \notin \text{gr}(L)$ . Then  $\text{gr}(L) = \text{gr}_1(L) - \{\text{id}_1\}$ . Let  $\text{gr}(L) = h(M)$  for some interpretation  $h$  and some  $M \in K$ ,  $M \subseteq A^*$ . Let  $B$  be the set of all  $a \in A$  such that  $h(a)$  has at least one edge. Then the regular language  $A^*BA^*$  is the set of all  $w \in A^*$  such that  $h(w)$  contains at least one edge. Hence  $\text{gr}(L) = h(M \cap A^*BA^*) \in \text{Int}(K)$ .

Second, let  $\text{gr}_2(w) = \text{backfold}(\text{gr}_1(w))$ , and let  $\text{Str}_2$  be defined on the basis of  $\text{gr}_2$ . Then  $\text{Str}_2(\text{Int}(K)) = \text{Str}_1(\text{Int}(K))$ . This is because for every graph language  $L$ ,  $L \in \text{Int}(K)$  iff  $\text{backfold}(L) \in \text{Int}(K)$ . To see this, note that for every graph  $g$  with  $\text{type}(g) = (m, n)$ ,  $\text{backfold}(g) = (g \oplus \text{id}_n) \circ \text{backfold}(\text{id}_n)$  (see the proof of Theorem 7) and  $g = (\text{id}_m \oplus \text{fold}(\text{id}_n)) \circ (\text{backfold}(g) \oplus \text{id}_n)$  (see the proof of Theorem 13). Thus, it suffices to show that  $\text{Int}(K)$  is closed under the operations  $L' \oplus \{\text{id}_n\}$ ,  $\{h\} \circ L'$ , and  $L' \circ \{h\}$ . For the first operation this has been shown in the proof of Lemma 19. The other two operations are left to the reader (see the end of the proof of Theorem 13 and the end of the proof of Theorem 27).

Third, define  $\text{gr}_3(w)$  in the same way as  $\text{gr}_2(w)$ , except that the type of the edges is changed from  $(1, 1)$  to  $(2, 0)$ . To be precise, an edge  $e$  with  $s(e) = u$  and  $t(e) = v$  in  $\text{gr}_2(w)$ , has  $s(e) = uv$  and  $t(e) = \lambda$  in  $\text{gr}_3(w)$ . It should be clear that  $\text{Str}_3(\text{Int}(K)) = \text{Str}_2(\text{Int}(K))$ , where  $\text{Str}_3$  is based on  $\text{gr}_3$  in the usual way.

In [EngHey], the coding  $\text{gr}_3$  is used instead of gr. We have just shown that this does not change the class  $\text{Str}(\text{LIN-HR}(K))$ , in the sense that  $\text{Str}_3(\text{LIN-HR}(K)) =$

$\text{Str}(\text{LIN-HR}(K))$ . Another small difference is that in [EngHey] all (terminal and nonterminal) edges of graphs have type  $(m, 0)$  for some  $m$ . Let us indicate this here by  $\text{HR}'$ . Since all edges of  $\text{gr}_3(w)$  have type  $(2, 0)$ , and since  $\text{backfold}(\text{gr}_3(w)) = \text{gr}_3(w)$ , it should be clear from the construction of  $\overline{G}$  in the proof of Theorem 13 that for every language  $L$ ,  $\text{gr}_3(L) \in \text{LIN-HR}'(K)$  iff  $\text{gr}_3(L) \in \text{LIN-HR}(K)$ . This shows that  $\text{Str}(\text{LIN-HR}(K)) = \text{Str}_3(\text{LIN-HR}'(K))$ , the class considered in [EngHey].

It is proved in [EngHey] that  $\text{Str}(\text{LIN-HR})$  equals the class of ranges of 2dgsms mappings. Since it is easy to see that  $\text{LIN-HR}(\text{REG}) = \text{LIN-HR}$  and that  $2\text{DGSM}(\text{REG})$  is the class of output languages of 2dgsms's (by incorporating the regular control language in the finite control of the grammar and the 2dgsms, respectively), this proves the theorem for  $K = \text{REG}$ . The proof of the general case consists of a careful analysis of the proof in [EngHey], which shows that it can be generalized to  $K$ -controlled grammars, under the assumption that  $K$  is closed under sequential machine mappings (and  $K$  is nontrivial). This analysis can easily be carried out.  $\square$

From Theorems 27, 28, and Lemma 29, we obtain our second characterization of the class  $\text{Ext}(K)$ .

**Theorem 30.** *For every nontrivial class  $K$  that is closed under sequential machine mappings,  $\text{Ext}(K) = 2\text{DGSM}(K)$ .*

**Corollary 31.** *For all nontrivial classes  $K$  and  $K'$  that are closed under sequential machine mappings,  $\text{Int}(K) = \text{Int}(K')$  if and only if  $2\text{DGSM}(K) = 2\text{DGSM}(K')$ .*

Quite a lot is known about the class  $2\text{DGSM}(K)$ , see, e.g., [EngRS]. As an example, it equals the class of languages generated by  $K$ -controlled ETOL systems of finite index (Corollary 4.10 of [EngRS]). The trivial fact that  $\text{Ext}(\text{Ext}(K)) = \text{Ext}(K)$  corresponds to the known result that  $2\text{DGSM}(2\text{DGSM}(K))$  is equal to  $2\text{DGSM}(K)$ ; this shows that  $\text{Ext}(K)$  is closed under 2dgsms mappings (cf. Corollary 5.8 of [EngRS]).

Theorem 30 and Corollary 31 allow us to use known formal language theoretic results for the classes  $2\text{DGSM}(K)$  to find out the power of the classes  $\text{Int}(K)$ . Thus, for  $K = \text{REG}$ ,  $\text{Ext}(K)$  is the class  $2\text{DGSM}(\text{REG})$  of output languages of 2dgsms mappings. Since it is well known that the class DB of derivation-bounded context-free languages is contained in  $2\text{DGSM}(\text{REG})$  (see, e.g., [Raj]), this implies the previously mentioned result that  $\text{Int}(\text{DB}) = \text{Int}(\text{REG})$ . Also, since there is a context-free language not in  $2\text{DGSM}(\text{REG})$ , see Lemma 4.24 of [Gre] (or Theorem 3.2.17 of [EngRS]),  $\text{Int}(\text{REG})$  is properly included in  $\text{Int}(\text{CF})$ .

**Theorem 32.**  $\text{Int}(\text{REG}) = \text{Int}(\text{LIN}) = \text{Int}(\text{DB}) \subset \text{Int}(\text{CF}) \subset \text{HR}$ .

Finally, we would like to know whether  $\text{Int}(\text{CF})$  is the largest class  $\text{Int}(K)$  that is included in HR. This is true if and only if  $\text{Ext}(\text{CF})$  is the largest class  $K$  such that  $\text{Int}(K)$  is included in HR. Trying to find an answer to this question, we first characterize this largest class.

**Theorem 33.**  $\text{Str}(\text{HR})$  is the largest class  $K$  such that  $\text{Int}(K) \subseteq \text{HR}$ .

*Proof.* The proof is similar to the one of Theorem 28. We first observe that HR is closed under replacements. In fact, if  $G$  is a context-free graph grammar over CS (see Theorem 11) and  $\phi$  is a replacement, then  $\phi(\text{val}(L(G))) = \text{val}(L(G'))$ , where  $G'$  is obtained from  $G$  by changing every constant  $c_g$  that occurs in the productions of  $G$  into  $c_{\phi(g)}$ . The correctness of this construction follows from Lemma 4.

As in the proof of Theorem 28 it now suffices to show that  $\text{AtInt}(\text{Str}(\text{HR})) \subseteq \text{HR}$ . Instead of HR we consider the class  $\text{Val}(\text{CFG}(\text{CS}))$ , see Theorem 11. Let  $G = (N, T, P, S)$  be a context-free graph grammar over CS such that  $\text{val}(L(G)) = \text{gr}(L)$  for some language  $L \subseteq B^*$ , and let  $h$  be an atomic interpretation of  $B$  such that  $h(L)$  is a graph language. Note that, as in the proof of Theorem 28, each symbol  $b \in B$  has type  $(1, 1)$  in  $\text{val}(L(G))$ , and has an arbitrary type with respect to  $h$ . We may assume that  $G$  is in *normal form*, i.e., that all its productions are of the form  $X \rightarrow c_g$  or  $X \rightarrow Y \circ Z$  or  $X \rightarrow Y \oplus Z$ , where  $X, Y, Z \in N$  and  $g \in \text{GR}$  (this is the usual normal form of regular tree grammars, see, e.g., [GécSte]).

We have to construct a context-free graph grammar  $G' = (N', T', P', S')$  such that  $\text{val}(L(G')) = h(L)$ . The idea of the construction is the same as in the proof of Theorem 28: each graph  $\text{gr}(w)$  for which  $h(w)$  is defined, is transformed into  $\text{exp}(\text{gr}(w), \mu)$ , where  $\mu$  is the unique decoration of  $\text{gr}(w)$  (see the proof of Theorem 28 for the terminology used). Let  $M$  be the maximal number occurring in the types of the symbols of  $B$  (with respect to  $h$ ). We define  $N'$  to consist of all triples  $(X, \mu_b, \mu_e)$  such that  $X \in N$  and  $\mu_b, \mu_e \in \{1, \dots, M\}^*$  with  $\text{type}(X) = (|\mu_b|, |\mu_e|)$ ; moreover,  $\text{type}(X, \mu_b, \mu_e) = \text{type}(X)$ . The intuition is that if  $\text{atom}(X) \Rightarrow^* g$  where  $g$  is terminal, then  $\text{atom}(X, \mu_b, \mu_e) \Rightarrow^* \text{exp}(g, \mu)$  where  $\mu$  is the decoration of  $g$  such that  $\mu(\text{begin}(g)(i)) = \mu_b(i)$  for all  $1 \leq i \leq |\mu_b|$  and  $\mu(\text{end}(g)(j)) = \mu_e(j)$  for all  $1 \leq j \leq |\mu_e|$  (note that, assuming  $G$  to be reduced,  $\mu$  is unique because all isolated nodes of  $g$  are begin or end nodes). The initial nonterminal  $S'$  of  $G'$  is  $(S, \langle m \rangle, \langle n \rangle)$ , where  $(m, n) = \text{type}(h(L))$ . The productions in  $P'$  are defined as follows. If  $X \rightarrow Y \circ Z$  is in  $P$ , then  $(X, \mu_b, \mu_e) \rightarrow (Y, \mu_b, \mu) \circ (Z, \mu, \mu_e)$  is in  $P'$  for all appropriate strings  $\mu_b, \mu_e$ , and  $\mu$  over  $\{1, \dots, M\}$ . If  $X \rightarrow Y \oplus Z$  is in  $P$ , then  $(X, \mu_b \cdot \mu'_b, \mu_e \cdot \mu'_e) \rightarrow (Y, \mu_b, \mu_e) \oplus (Z, \mu'_b, \mu'_e)$  is in  $P'$  for all  $(Y, \mu_b, \mu_e), (Z, \mu'_b, \mu'_e) \in N'$ . Finally, if  $X \rightarrow c_g$  is in  $P$ , then  $(X, \mu_b, \mu_e) \rightarrow c_{\text{exp}(g, \mu)}$  is in  $P'$  for all strings  $\mu_b, \mu_e$  and all decorations  $\mu$  of  $g$  such that  $\mu(\text{begin}(g)(i)) = \mu_b(i)$  and  $\mu(\text{end}(g)(j)) = \mu_e(j)$  for all appropriate  $i$  and  $j$ .

This ends the construction of  $G'$ . A formal correctness proof is left to the reader. It should be based on a definition of the sum of decorated graphs, and the fact that ‘exp’ is a homomorphism with respect to this sum (and the corresponding fact for the concatenation of decorated graphs, as observed in the proof of Theorem 28).  $\square$

This proves that  $\text{Int}(\text{Str}(\text{HR}))$  is the largest class  $\text{Int}(K)$  that is included in HR.

It is shown in [EngHey] that the class  $\text{Str}(\text{HR})$  of string languages generated by HR grammars is equal to the class  $\text{OUT}(\text{DTWT})$  of output lan-



guages of deterministic tree-walking transducers. It now follows from Theorems 30 and 33 that  $\text{Int}(\text{CF})$  is the largest class  $\text{Int}(K)$  that is included in HR if and only if  $\text{Ext}(\text{CF})$  is the largest class  $K$  such that  $\text{Int}(K) \subseteq \text{HR}$  if and only if  $2\text{DGSM}(\text{CF}) = \text{OUT}(\text{DTWT})$ . Note that it follows from our results that  $2\text{DGSM}(\text{CF}) = \text{Ext}(\text{CF}) = \text{Str}(\text{Int}(\text{CF})) \subseteq \text{Str}(\text{HR}) = \text{OUT}(\text{DTWT})$ , which was proved in a completely different way in Corollary 5.6 of [EngRS] (where  $2\text{DGSM}(\text{CF})$  is denoted  $\text{DCS}(\text{CF})$ , and  $\text{OUT}(\text{DTWT})$  is denoted  $\text{DCT}(\text{REC})$  or  $yT_{\text{fc}}(\text{REC})$ ). However, equality of  $2\text{DGSM}(\text{CF})$  and  $\text{OUT}(\text{DTWT})$  is mentioned as an open problem after Corollary 5.6 of [EngRS]. Hence, it is an *open problem* whether or not  $\text{Int}(\text{CF})$  is the largest class  $\text{Int}(K)$  that is included in HR. As *another open problem* we mention the following: is it true that every HR graph language of bounded pathwidth is in  $\text{Int}(\text{Str}(\text{HR}))$ ? Or even in  $\text{Int}(\text{CF})$ ? Note that all HR graph languages in  $\text{Int}(\text{Str}(\text{HR}))$  are of bounded pathwidth by Corollary 17. Some *more open problems* are: is it decidable whether an HR graph language is in  $\text{Int}(\text{REG})$ ? and the same question for  $\text{Int}(\text{CF})$  and  $\text{Int}(\text{Str}(\text{HR}))$ .

We finally mention that it would be interesting to find another natural operation of concatenation of graphs that can be used to characterize the graph languages generated by the linear edNCE grammars (which are node replacement graph grammars, see, e.g., [CouER]).

**Acknowledgment.** We wish to thank Hans Bodlaender for the references to pathwidth.

## References

- [BauCou] M.Bauderon, B.Courcelle; Graph expressions and graph rewritings, *Math. Syst. Theory* 20 (1987), 83-127
- [Ben] D.B.Benson; The basic algebraic structures in categories of derivations, *Inform. and Control* 28 (1975), 1-29
- [Ber] J.Berstel; *Transductions and Context-Free Languages*, Teubner, Stuttgart, 1979
- [Bod] H.L.Bodlaender; A partial  $k$ -arboretum of graphs with bounded treewidth, Preliminary version, Utrecht University, September 1995
- [BosDW] F.Bossut, M.Dauchet, B.Warin; A Kleene theorem for a class of planar acyclic graphs, *Inform. and Comput.* 117 (1995), 251-265
- [Cla] V.Claus; Ein Vollständigkeitssatz für Programme und Schaltkreise, *Acta Informatica* 1 (1971), 64-78
- [Cou1] B.Courcelle; An axiomatic definition of context-free rewriting and its application to NLC graph grammars, *Theor. Comput. Sci.* 55 (1987), 141-181
- [Cou2] B.Courcelle; Graph rewriting: an algebraic and logic approach, in *Handbook of Theoretical Computer Science, Vol.B* (J.van Leeuwen, ed.), Elsevier, 1990, pp.193-242
- [Cou3] B.Courcelle; The monadic second-order logic of graphs III: Tree-decompositions, minors and complexity issues, *RAIRO Theoretical Informatics and Applications* 26 (1992), 257-286
- [CouER] B.Courcelle, J.Engelfriet, G.Rozenberg; Handle-rewriting hypergraph languages, *J. of Comp. Syst. Sci.* 46 (1993), 218-270

- [Dre] F.Drewes; Transducibility - symbolic computation by tree-transductions, University of Bremen, Bericht Nr. 2/93, 1993
- [EhrKKK] H.Ehrig, K.-D.Kiermeier, H.-J.Kreowski, W.Kühnel; *Universal Theory of Automata*, Teubner, Stuttgart, 1974
- [EllST] J.A.Ellis, I.H.Sudborough, J.S.Turner; The vertex separation and search number of a graph, *Inform. and Comput.* 113 (1994), 50-79
- [Eng] J.Engelfriet; Graph grammars and tree transducers, Proc. CAAP'94 (S.Tison, ed.), *Lecture Notes in Computer Science 787*, Springer-Verlag, Berlin, 1994, pp.15-36
- [EngHey] J.Engelfriet, L.M.Heyker; The string generating power of context-free hypergraph grammars, *J. of Comp. Syst. Sci.* 43 (1991), 328-360
- [EngRS] J.Engelfriet, G.Rozenberg, G.Slutzki; Tree transducers, L systems, and two-way machines, *J. of Comp. Syst. Sci.* 20 (1980), 150-202
- [EngVer] J.Engelfriet, J.J.Vereijken; Concatenation of graphs, in *Graph-Grammars and their Application to Computer Science*, Proceedings of the 5th International Workshop, Williamsburg, 1994, to appear as *Lecture Notes in Computer Science*, Springer-Verlag, Berlin
- [GécSte] F.Gécseg, M.Steinby; *Tree Automata*, Akadémiai Kiadó, Budapest, 1984
- [Gre] S.Greibach; One-way finite visit automata, *Theor. Comput. Sci.* 6 (1978), 175-221
- [Hab] A.Habel; *Hyperedge Replacement: Grammars and Languages*, Lecture Notes in Computer Science 643, Springer-Verlag, Berlin, 1992
- [HabKre] A.Habel, H.-J.Kreowski; May we introduce to you: hyperedge replacement, in *Graph-Grammars and their Application to Computer Science* (H.Ehrig, M.Nagl, G.Rozenberg, A.Rosenfeld, eds.), *Lecture Notes in Computer Science 291*, Springer-Verlag, Berlin, 1987, pp.15-26
- [HabKV] A.Habel, H.-J.Kreowski, W.Vogler; Metatheorems for decision problems on hyperedge replacement graph languages, *Acta Informatica* 26 (1989), 657-677
- [HopUll] J.E.Hopcroft, J.D.Ullman; *Introduction to Automata Theory, Languages, and Computation*, Addison-Wesley, Reading, Mass., 1979
- [Hot1] G.Hotz; Eine Algebraisierung des Syntheseproblems von Schaltkreisen, *EIK* 1 (1965), 185-205, 209-231
- [Hot2] G.Hotz; Eindeutigkeit und Mehrdeutigkeit formaler Sprachen, *EIK* 2 (1966), 235-246
- [HotKM] G.Hotz, R.Kolla, P.Molitor; On network algebras and recursive equations, in *Graph-Grammars and Their Application to Computer Science* (H.Ehrig, M.Nagl, G.Rozenberg, A.Rosenfeld, eds.), *Lecture Notes in Computer Science 291*, Springer-Verlag, Berlin, 1987, pp.250-261
- [Klo] T.Kloks; *Treewidth*, Lecture Notes in Computer Science 842, Springer-Verlag, Berlin, 1994
- [Lau] C.Lautemann; Decomposition trees: structured graph representation and efficient algorithms, Proc. CAAP'88, *Lecture Notes in Computer Science 299*, Springer-Verlag, Berlin, 1988, pp.28-39
- [MezWri] J.Mezei, J.B.Wright; Algebraic automata and context-free sets, *Inform. and Control* 11 (1967), 3-29
- [Raj] V.Rajlich; Absolutely parallel grammars and two-way finite state transducers, *J. of Comp. Syst. Sci.* 6 (1972), 324-342
- [RobSey] N.Robertson, P.D.Seymour; Graph minors I. Excluding a forest, *J. Comb. Theory Ser.B* 35 (1983), 39-61
- [Sal] A.Salomaa; *Formal Languages*, Academic Press, New York, 1973

[Ver] J.J.Vereijken; *Graph Grammars and Operations on Graphs*, Master's Thesis,  
Leiden University, May 1993

This article was processed using the L<sup>A</sup>T<sub>E</sub>X macro package with LLNCS style