# No Rule is Forever: Datalog Reasoning with Rule Amendments (proof appendix)

Weiqin Xu
*LIGM Univ Gustave Eiffel, CNRS*
Marne la Vallée, France
weiqin.xu@univ-eiffel.fr

Riccardo Tommasini
*LIRIS, INSA*
Lyon, France
riccardo.tommasini@insa-lyon.fr

Olivier Curé
*LIGM Univ Gustave Eiffel, CNRS*
Marne la Vallée, France
olivier.cure@univ-eiffel.fr

## REFERENCES

[1] S. Abiteboul, R. Hull, and V. Vianu. *Foundations of Databases*. Addison-Wesley, 1995.

## APPENDIX

This appendix provides detailed proofs of the various lemmas of this article. For the reader's convenience, the lemmas are restated. We begin by providing some additional definitions and results adapted from [1].

### A. Additional definitions and results

For a semi-positive Datalog program $P$ and an extensional database EDB over $\mathrm{edb}(P)$, we define $P(\mathrm{EDB})$ as the set of facts derivable from EDB using the rules of $P$. This semantics is defined as the smallest fixed point of the direct consequence relation $T_P$. For any instance $K$, $T_P(K)$ contains the facts present in $K$ for predicates in $\mathrm{edb}(P)$, as well as the facts that can be derived from $P$ using the facts in $K$. For an extensional database EDB over $\mathrm{edb}(P)$, $P(\mathrm{EDB})$ corresponds to the limit of the increasing sequence $(T_P^k(\mathrm{EDB}))_{k \geq 0}$.

For convenience, we extend the definition of the semantics of $P$ to any instance $I$ by taking $P(I)$ to be $P(I \mid \mathrm{edb}(P))$, where $I \mid \mathrm{edb}(P)$ denotes the restriction of the instance $I$ to the relations in $\mathrm{edb}(P)$.

The following lemma recalls some basic properties of the direct consequence relation $T_P$:

**Lemma A.1.** *Let $P$ and $P'$ be two semi-positive Datalog programs such that $P' \subseteq P$ and $\mathrm{edb}(P') \subseteq \mathrm{edb}(P)$.*
1) *For any instance $I$, $T_P(I)|\mathrm{edb}(P) = I|\mathrm{edb}(P)$.*
2) *For any instance $I$ over $\mathrm{edb}(P)$, $P(I)|\mathrm{edb}(P) = I$.*
3) *For two instances $J' \subseteq J$ such that $J|\mathrm{edb}(P') = J'|\mathrm{edb}(P')$, we have $T_{P'}(J') \subseteq T_P(J)$.*

The following lemma relates the semantics of a program to that of its subprograms.

**Lemma A.2.** *Let $P' \subseteq P$ be semi-positive Datalog programs and let EDB be an extensional database over $\mathrm{edb}(P)$. We have*

$$P'(\mathrm{EDB}) \subseteq P(\mathrm{EDB}).$$

*Proof.* We prove by induction that for all $k \geq 0$, $T_{P'}^k(\mathrm{EDB}) \subseteq T_P^k(\mathrm{EDB})$. The base case is immediate. For the induction step, assume that $T_{P'}^k(\mathrm{EDB}) \subseteq T_P^k(\mathrm{EDB})$. Using Lemma A.1.1, we can easily show that

$$\begin{aligned} &T_{P'}^k(\mathrm{EDB}) \mid \mathrm{edb}(P') \\ = \ &T_P^k(\mathrm{EDB}) \mid \mathrm{edb}(P') \\ = \ &\mathrm{EDB} \mid \mathrm{edb}(P'). \end{aligned}$$

Applying Lemma A.1.3, we obtain $T_{P'}^{k+1}(\mathrm{EDB}) \subseteq T_P^{k+1}(\mathrm{EDB})$. It follows that $P'(\mathrm{EDB})$, which is the limit of the sequence $(T_{P'}^k(\mathrm{EDB}))_{k \geq 0}$, is included in $P(\mathrm{EDB})$, which is the limit of the sequence $(T_P^k(\mathrm{EDB}))_{k \geq 0}$. $\square$

### B. Proof of Lemma 1

**Lemma 1.** *A Hyper Rule Dependency Graph will always be a Directed Acyclic Graph (DAG).*

*Proof.* We define the expression $a \mapsto b$, which represents that there exists a path from $a$ to $b$ in a directed graph. Assuming an $HRDG$ is not a DAG, then there must exist a limited number of hypernodes $h1, h2, .., hk \in H$ ($k > 1$), where $h1 \mapsto h2 \mapsto ... \mapsto hk \mapsto h1$. According to the definition of $HRDG$, considering the positive case, if there exists a positive dependency between $h1$ and $h2$, then there must be two rules $r_{12} \in h1$ and $r_{21} \in h2$ and $(r_{12}, r_{21}) \in P$ ($P$ represents the positive $RDG$), where $r_{12} \mapsto r_{21}$. Similarly, if $h2 \mapsto h3$, there must be two rules $r_{22} \in h2$ and $r_{31} \in h3$ and $(r_{22}, r_{31}) \in P$, where $r_{22} \mapsto r_{31}$. Moreover, according to the definition of a hypernode, we either have $r_{21} \mapsto r_{22}$ or $r_{21} = r_{22}$. Thus, if $h1 \mapsto h2 \mapsto ... \mapsto hk \mapsto h1$, then there must exist $r_{11}, r_{12} \in h1, r_{21}, r_{22} \in h2, ..., r_{k1}, r_{k2} \in hk$ where $r_{11} \mapsto r_{12} \mapsto r_{21} \mapsto r_{22} \mapsto ... \mapsto r_{k1} \mapsto r_{k2} \mapsto r_{11}$, meaning $r_{11}, r_{12}, r_{21}, r_{22}, ..., r_{k1}, r_{k2}$ belong to a single hypernode instead of different ones, which is contradictory to the original assumption. Similarly, for the strictly negative and positive/negative cases, an $HRDG$ is a DAG. $\square$

### C. Proof of correctness of Algorithm 1

The aim of this section is to show that Algorithm 1 correctly computes the semantics of a stratified Datalog program.

We begin by recalling the definition of the semantics of a stratified Datalog program. Let $P$ be a stratified Datalog program and let $\sigma$ be one of its stratifications which induces a decomposition $P_1, \ldots, P_n$ of the program $P$. The semantics of $P$ under the stratification $\sigma$ with an extensional database EDB,

denoted $P^\sigma(\text{EDB})$, is defined as the limit of the increasing sequence $(I_k)_{k\in[0,n]}$ defined by:

$$
\begin{aligned}
I_0 &= \text{EDB} \\
I_i &= I_{i-1} \cup P_i(I_{i-1} \mid \text{edb}(P_i)) \qquad \text{for } i \in [1, n]
\end{aligned}
$$

The semantics $P^{\text{strat}}(\text{EDB})$ of a stratified Datalog program $P$ over an extensional database EDB is defined as $P^\sigma(\text{EDB})$ for any stratification $\sigma$ of $P$ (see [1][Lemma 15.2.8]).

**Lemma A.3.** *For a stratified Datalog program $P$ and an extensional database* EDB *over* $\text{edb}(P)$*, Algorithm 1 computes* $P^{\text{strat}}(\text{EDB})$*.*

*Proof.* Let $P$ be a stratified Datalog program and let HRDG be its associated hypergraph. Let $\text{hn}_1, \ldots, \text{hn}_n$ be an enumeration of the hypernodes of HRDG following the topological ordering computed at the beginning of Algorithm 1 Algorithm 1 computes the sequence $(\text{IDB}_{\text{hn}_k})_{k\in[1,n]}$ defined by:

$$
\text{IDB}_{\text{hn}_k} = \text{hn}_k(\text{EDB}_{\text{hn}_k} \mid \text{edb}(\text{hn}_k))
$$

where $\text{EDB}_{\text{hn}_k} = \bigcup_{(\text{hn}',\text{hn}_k)\in\text{HRDG}} \text{IDB}_{\text{hn}'} \cup \text{EDB}$. The value computed by Algorithm 1 is the union $\bigcup_{k\in[1,n]} \text{IDB}_{\text{hn}_k}$.

Let us first remark that the topological ordering $\text{hn}_1, \ldots, \text{hn}_n$ of the hypernodes of HRDG naturally induces a stratification $\sigma$ of $P$. For a rule $r$ of $P$, the stratification $\sigma(r)$ is defined to be the unique $k \in [1, n]$ such that $r$ belongs to $\text{hn}_k$. Assume towards a contradiction that $\sigma$ is not a stratification of $P$. There are two possibilities:

- Either there would exists a positive edge $(r, r')$ with $r \in \text{hn}_i$ and $r' \in \text{hn}_j$ with $i > j$ which contradicts the definition of a topological ordering.
- Or there would exists a negative edge $(r, r')$ with $r \in \text{hn}_i$ and $r' \in \text{hn}_j$ with $i \geq j$. As we started with a topological ordering, it must be the case that $i = j$. As the edge is negative, this would imply that $\text{hn}_i$ contains a cycle with a negative edge which contradicts the fact that $P$ can be stratified.

The semantics of $P$ under the stratification $\sigma$ is the limit of the increasing sequence $(I_k)_{k\in[0,n]}$ defined by:

$$
\begin{aligned}
I_0 &= \text{EDB} \\
I_i &= I_{i-1} \cup P_i(I_{i-1} \mid \text{edb}(P_i)) \qquad \text{for } i \in [1, n]
\end{aligned}
$$

We will prove by induction on $k \in [1, n]$ that:

$$
I_k = \bigcup_{i\in[1,k]} \text{IDB}_{\text{hn}_i} \cup \text{EDB}
$$

For $k = 1$, the result is immediate. For the induction step, assume that $I_k = \bigcup_{i\in[1,k]} \text{IDB}_{\text{hn}_i} \cup \text{EDB}$ for some $k \geq 1$. To show that $I_{k+1} = \bigcup_{i\in[1,k+1]} \text{IDB}_{\text{hn}_i} \cup \text{EDB}$, it is enough to show that

$$
\text{hn}(I_k \mid \text{edb}(\text{hn}_k)) = \text{hn}(\text{EDB}_{\text{hn}_k} \mid \text{edb}(\text{hn}_k)).
$$

Clearly $\text{EDB}_{\text{hn}_k} \subseteq I_k$. Now consider a fact in $I_k$ for a predicate in $\text{edb}(\text{hn}_k)$. This fact either belongs to EDB or, by definition of HRDG to some $\text{hn}'$ with $(\text{hn}', \text{hn}_k) \in \text{HRDG}$.

By definition of $\text{EDB}_{\text{hn}_k}$, it follows that $I_k \mid \text{edb}(\text{hn}_k) = \text{EDB}_{\text{hn}_k} \mid \text{edb}(\text{hn}_k)$. This concludes the induction proof. We have established that:

$$
P^{\text{strat}}(\text{EDB}) = P^\sigma(\text{EDB}) = \bigcup_{i\in[1,k]} \text{IDB}_{\text{hn}_i} \cup \text{EDB}.
$$

To conclude the proof, it suffices to remark that $\text{EDB} \subseteq \bigcup_{k\in[1,n]} \text{IDB}_{\text{hn}_k}$. Indeed $\text{EDB} = \bigcup_{k\in[1,n]} \text{EDB} \mid \text{edb}(\text{hn}_k)$ and for all $k \in [1, n]$, $\text{EDB} \mid \text{edb}(\text{hn}_k) \subseteq \text{hn}_k(\text{EDB}_{\text{hn}_k})$ by Lemma A.1.2 as $\text{EDB} \subseteq \text{EDB}_{\text{hn}_k}$. $\qquad\square$

*D. Proof of Lemma 2*

**Lemma 2.** *Let $P$ be a semi-positive Datalog program and an extensional database $EDB$ over* $\text{edb}(P)$*. For $B$ a subset of* $P(\text{EDB})$*,*

$$
P(\text{EDB}) = \text{SemiNaiveAlgorithm}(P, \text{EDB} \cup B)
$$

*Proof.* To establish this equality, we prove that the limit of the increasing sequence $(T_P^k(\text{EDB}\cup B))_{k\geq 1}$ is equal to $P(\text{EDB})$.

We will show by induction on $k \geq 0$ that:

$$
T_P^k(\text{EDB}) \subseteq T_P^k(\text{EDB} \cup B) \subseteq T_P^k(P(\text{EDB})).
$$

For $k = 0$, we do have that $\text{EDB} \subseteq \text{EDB} \cup B \subseteq P(\text{EDB})$. For the induction step, assume that $T_P^k(\text{EDB}) \subseteq T_P^k(\text{EDB}\cup B) \subseteq T_P^k(P(\text{EDB}))$. By Lemma A.1.1, we have that:

$$
\begin{aligned}
T_P^k(\text{EDB}) \mid \text{edb}(P) &= \text{EDB} \mid \text{edb}(P) \\
T_P^k(P(\text{EDB})) \mid \text{edb}(P) &= P(\text{EDB}) \mid \text{edb}(P) \\
&= \text{EDB} \mid \text{edb}(P) \\
T_P^k(\text{EDB} \cup B) \mid \text{edb}(P) &= \text{EDB} \cup B \mid \text{edb}(P)
\end{aligned}
$$

As $\text{EDB} \subseteq \text{EDB} \cup B \subseteq P(\text{EDB})$, it follows that $\text{EDB} \cup B \mid \text{edb}(P) = \text{EDB} \mid \text{edb}(P)$. Hence, we have:

$$
\begin{aligned}
&T_P^k(\text{EDB}) \mid \text{edb}(P) \\
=\ &T_P^k(\text{EDB} \cup B) \mid \text{edb}(P) \\
=\ &T_P^k(P(\text{EDB})) \mid \text{edb}(P) \\
=\ &\text{EDB} \mid \text{edb}(P)
\end{aligned}
$$

Hence we can apply Lemma A.1.3 to conclude the induction.

As $(T_P^k(\text{EDB}))_{k\geq 0}$ and $T_P^k(P(\text{EDB}))$ share the same limit $P(\text{EDB})$, we have that the limit of $(T_P^k(\text{EDB} \cup B))_{k\geq 0}$ is equal to $P(\text{EDB})$. $\qquad\square$

*E. Proof of correctness in the case of rule insertions*

We will now show that the semantics computed when adding rules is correct.

We start by recalling the procedure in the case of rule insertion. We start with a semi-positive stratified Datalog program $P$ with its HRDG $G$. Furthermore, the semantics $\text{IDB}_{\text{hn}}$ for $P$ of each hypernode hn of $G$ is assumed to have been computed.

We consider a program $P^t$ obtained by adding a new set of rules $NR$ to $P$. We compute its HRDG $G^t$ and the semantics $\text{IDB}^t_{hn}$ of each of its hypernode hn as follows:

- using Algorithm 2, we compute $G^t$ together with a set DIH of directly impacted hypernodes of $G^t$. For each hypernode $\text{hn} \in \text{DIH}$, Algorithm 2 provides a

bootstrapping value noted here $\text{Boot}_{\text{hn}}$ that will be used to speed up the computation of $\text{IDB}_{\text{hn}}^t$.

- using Algorithm 4, we compute a topological ordering on the set IH of hypernodes of $G^t$ reachable for DIH.
- the semantics $\text{IDB}_{hn}^t$ for hypernodes hn in $G^t$ are computed as follows:
  - First, for hypernodes hn that do not belong to DIH $\cup$ IH, $\text{IDB}_{\text{hn}}^t$ is taken to be equal to $\text{IDB}_{\text{hn}}$.
  - The values of $\text{IDB}_{\text{hn}}^t$ for hn $\in$ DIH$\cup$IH are computed using the semi-naive algorithm using the topological ordering previously constructed as in Algorithm 1. For hypernodes hn in DIH the computation of the semi-naive algorithm is initialized using $\text{Boot}_{\text{hn}}$ instead of a empty set in the other case.

For an HRDG $G$ and one of its hypernodes hn, we denote by $G{\downarrow}_{\text{hn}}$ the HRDG obtained by restricting $G$ to the hypernodes that can reach hn in $G$. The correctness of our procedure relies on the following key property:

**Lemma A.4.** *For a hypernode* hn *that does not belong to* DIH $\cup$ IH,

$$G^t{\downarrow}_{hn}= G{\downarrow}_{\text{hn}} \ .$$

*Proof.* As $P'$ is obtained by adding rules from NR to $P$, the hypernodes of $G^t$ are only of two form:

- a hypernode of $G$
- the union of a non-empty subset of NR with a possibly empty union of hypernodes of $G$.

Furthermore, the vertices in DIH are precisely the hypernodes of $G^t$ containing a rule in NR.

Now consider a hypernode hn that does not belong to IH $\cup$ DIH. As hn does not belong to DIH, hn does not contain a rule in NR and from what precedes hn belongs to $G$. As hn does not belong to IH, it is not reachable from any hypernode containing a rule in NR. Hence all hypernodes in $G^t{\downarrow}_{\text{hn}}$ belong to $G$. It follows that $G{\downarrow}_{\text{hn}}{\subseteq} G^t{\downarrow}_{\text{hn}}$. Assume towards a contradiction that the inclusion is strict. That is to say that there exists a hypernode hn$'$ that can reach hn in $G$ but that does not belongs to $G^t$. This would imply that there exists a hypernode hn$''$ in $G^t$ that contain hn$'$ and at least one rule in NR. Hence hn$''$ would belong to DIH and as hn$''$ contains hn$'$, it can reach hn also in $G^t$ which contradicts the fact that all hypernodes that can reach hn in $G^t$ also belong to $G$. $\qquad\square$

Armed with Lemma A.4, we can show that our procedure is correct.

**Lemma A.5.** *The incremental update procedure is correct when adding rules.*

*Proof.* Clearly Algorithm 2 computes the HRDG $G^t$. Also Algorithm 4 is an adaptation to hypergraph of a the standard DFS based computation of a topological ordering. So we only need to show that the computation of the various $\text{IDB}_{\text{hn}}^t$ for hn $\in G^t$ is correct.

We consider hypernodes in a topological ordering that extends the ordering computed by Algorithm 4 and show by induction that the computed value for $\text{IDB}_{\text{hn}}^t$ is correct.

**If** hn **does not belong to** DIH $\cup$ IH**.** In this case, we take $\text{IDB}_{\text{hn}}^t = \text{IDB}_{\text{hn}}$. By Lemma A.4, we have that $G^t{\downarrow}_{\text{hn}}= G{\downarrow}_{\text{hn}}$. As $\text{IDB}_{\text{hn}}^t$ only depends on $G^t{\downarrow}_{\text{hn}}$, it follows that $\text{IDB}_{\text{hn}}^t$ is indeed equal to $\text{IDB}_{\text{hn}}$.

**If** hn **belongs to** DIH**.** We take $\text{EDB}_{\text{hn}}^t$ to be the result of the semi-naive algorithm applied to hn with:

$$\text{EDB}_{\text{hn}}^t \cup \text{Boot}_{\text{hn}}$$

as input instance where:

- $\text{EDB}_{\text{hn}}^t$ is the union of all $\text{EDB}_{\text{hn}'}^t$ for $(\text{hn}', \text{hn}) \in G^t$,
- $\text{Boot}_{\text{hn}}$ is computed in Algorithm 2.

As by induction hypothesis, all the $\text{EDB}_{\text{hn}'}^t$ for $(\text{hn}', \text{hn}) \in G^t$ are correct, we only need to show that

$$\text{SemiNaive}(hn, \text{EDB}_{\text{hn}}^t) = \text{SemiNaive}(hn, \text{EDB}_{\text{hn}}^t \cup \text{Boot}_{\text{hn}}).$$

By Lemma 2 , it is enough to show that $\text{Boot}_{\text{hn}}$ is included in $\text{hn}(\text{EDB}_{\text{hn}}^t)$. By construction, $\text{Boot}_{\text{hn}}$ is a finite union of $\text{IDB}_{\text{hn}_1}, \ldots, \text{IDB}_{\text{hn}_k}$ with for $i \in [1, k]$, $\text{hn}_i \subseteq \text{hn}$. We will show that for all $i \in [1, k]$, $\text{IDB}_{\text{hn}_i} \subseteq \text{IDB}_{\text{hn}}^t$. Let $i \in [1, k]$. We have that $\text{IDB}_{\text{hn}_i} = \text{hn}_i(\text{EDB}_{\text{hn}_i})$ and $\text{IDB}_{\text{hn}}^t = \text{hn}(\text{EDB}_{\text{hn}}^t)$. Using Lemma A.1.3, it is enough to shown that $\text{EDB}_{\text{hn}_i} \mid \text{edb}(\text{hn}_i) = \text{EDB}_{\text{hn}}^t \mid \text{edb}(\text{hn}_i)$. To see this, consider a relation $r \in \text{edb}(\text{hn}_i) \subseteq \text{edb}(\text{hn})$. Let hn be the hypernode in $G^t$ containing $r$. We have $(\text{hn}', \text{hn}) \in G^t$ and hence hn$'$ does not belong DIH $\cup$ IH. We have already shown that $G^t{\downarrow}_{\text{hn}'}= G{\downarrow}_{\text{hn}'}$ and that $\text{IDB}_{\text{hn}'} = \text{IDB}_{\text{hn}'}^t$. The facts corresponding to the relation $r$ in $\text{EDB}_{\text{hn}_i}$ are contained in $\text{IDB}_{\text{hn}'}$. Similarly the facts corresponding to the relation $r$ in $\text{EDB}_{\text{hn}}^t$ are contained in $\text{IDB}_{\text{hn}'}^t = \text{IDB}_{\text{hn}'}$.

**If** hn **belongs to** IH**.** $\text{EDB}_{\text{hn}}^t$ is computed exactly as in Algorithm 1 which has been shown to produce the correct semantics in Lemma A.3. $\qquad\square$

### F. Proof of correctness in the case of rule deletions

We will now show that the semantics computed when rules are removed is correct.

We start by recalling the procedure in the case of rule deletion. We start with a semi-positive stratified Datalog program $P$ with its HRDG $G$. Furthermore, the semantics $\text{IDB}_{\text{hn}}$ for $P$ of each hypernode hn of $G$ is assumed to have been computed.

We consider a program $P^t$ obtained by removing a set of rules $DR$ from $P$. We compute its HRDG $G^t$ and the semantics $\text{IDB}_{hn}^t$ of each of its hypernode hn as follows:

- using Algorithm 3 , we compute $G^t$ together with a set DIH of directly impacted hypernodes of $G^t$.
- using Algorithm 4 , we compute a topological ordering on the set IH of hypernodes of $G^t$ reachable for DIH.
- the semantics $\text{IDB}_{hn}^t$ for hypernodes hn in $G^t$ are computed as follows:
  - First, for hypernodes hn that do not belong to DIH $\cup$ IH, $\text{IDB}_{\text{hn}}^t$ is taken to be equal to $\text{IDB}_{\text{hn}}$.

– The values of $\text{IDB}^t_{\text{hn}}$ for $\text{hn} \in \text{DIH} \cup \text{IH}$ are computed using the semi-naive algorithm using the topological ordering previously constructed as in Algorithm 1.

The key property of Algorithm 3 is given by the following lemma.

**Lemma A.6.** *For a hypernode* $\text{hn}$ *that does not belong to* $\text{DIH} \cup \text{IH}$,

$$G^t{\downarrow}_{hn} = G{\downarrow}_{\text{hn}} .$$

*Proof.* As $P'$ is obtained by removing rules from $P$, the hypernodes of $G^t$ are only of two form:

- a hypernode of $G$,
- or strict subset of hypernodes of $G$ that contain rules in $DR$.

Furthermore, the vertices in DIH are precisely the hypernodes of $G^t$ that are strict subsets of hypernodes of $G$ and the hypernodes of $G$ that are successor in $G$ of hypernodes containing a rule in $DR$.

Now consider an hypernode $\text{hn}$ that does not belong to $\text{IH} \cup \text{DIH}$. As $\text{hn}$ does not belong to DIH, $\text{hn}$ belongs to $G$. It suffices to show that $G{\downarrow}_{\text{hn}}$ only contains hypernodes without rules in $DR$. Toward a contradiction, assume it is not the case. Let $\text{hn}'$ be a predecessor of $\text{hn}$ in $G$ containing a rule in $DR$ with the smallest distance to $\text{hn}$. By minimality of the distance, there exists a successor $\text{hn}''$ of $\text{hn}'$ which is a predecessor of $\text{hn}$ and which belong to $G$. By construction Algorithm 2, $\text{hn}''$ belongs to DIH and hence $\text{hn}$ belongs to $\text{DIH} \cup \text{IH}$ which

brings the contradiction. As $P'$ is obtained by removing rules from $P$, the desired equality follows from the fact that $G{\downarrow}_{\text{hn}}$ only contains hypernodes without rules in $DR$. $\square$

**Lemma A.7.** *The incremental update procedure is correct when deleting rules.*

*Proof.* Clearly Algorithm 3 computes the HRDG $G^t$. Also Algorithm 4 is an adaptation to hypergraph of a the standard DFS based computation of a topological ordering. So we only need to show that the computation of the various $\text{IDB}^t_{\text{hn}}$ for $\text{hn} \in G^t$ is correct.

We consider hypernodes in a topological ordering that extends the ordering computed by Algorithm 4 and show by induction that the computed value for $\text{IDB}^t_{\text{hn}}$ is correct.

**If** $\text{hn}$ **does not belong to** $\text{DIH} \cup \text{IH}$**.** In this case, we take $\text{IDB}^t_{\text{hn}} = \text{IDB}_{\text{hn}}$. By Lemma A.6, we have that $G^t{\downarrow}_{\text{hn}} = G{\downarrow}_{\text{hn}}$. As $\text{IDB}^t_{\text{hn}}$ only depends on $G^t{\downarrow}_{\text{hn}}$, it follows that $\text{IDB}^t_{\text{hn}}$ is indeed equal to $\text{IDB}_{\text{hn}}$.

**If** $\text{hn}$ **belongs to** $\text{DIH} \cup \text{IH}$**.** $\text{IDB}^t_{\text{hn}}$ is computed exactly as in Algorithm 1 which has been shown to produce the correct semantics in Lemma A.3. $\square$

*G. Proof of Lemma 3*

**Lemma 3.** *The incremental approach designed for rule amendments is correct.*

*Proof.* This is a direct consequence of Lemma A.5 and Lemma A.7. $\square$