

(Originally **WG3:W13-014**, “Definitions of path and subpath in SQL/PGQ and GQL”)  
**WG3:W16-045** pre-publication cross-post for WG3 document archive  
**LDBC Technical Report TR-2021-01**  
**Property graphs and paths in GQL: Mathematical definitions**

**WG3:W16-045** pre-publication cross-post for WG3 document archive

**LDBC Technical Report TR-2021-01**

## **Property graphs and paths in GQL: Mathematical definitions**

October 2021

Cite as:

Green, Alastair; Guagliardo, Paolo and Libkin, Leonid, 2021. “Property graphs and paths in GQL: Mathematical Definitions”, Linked Data Benchmark Council, Technical Report TR-2021-01. [\[DOI\]](#).

### **Abstract**

This paper provides precise mathematical definitions of a property graph as specified in the proposed GQL international standard, which is an attributed mixed multigraph with loops. It further defines a partially-oriented walk in such a property graph, which is called a path in GQL, as well as restricted classes of such walks (trails, simple/acyclic paths).

### **Note on publication as an LDBC Technical Report**

The original version of this document was submitted to ISO/IEC JTC1 SC32/WG3 (the committee which develops and maintains the SQL and GQL international standards) as document number **WG3:W13-014**, under the title “Definitions of path and subpath in SQL/PGQ and GQL”.

The original discussion paper, and this report, are both published on the Linked Data Benchmark Council (LDBC) website. Publication by LDBC does not imply that members of LDBC other than the authors necessarily agree with the content of this report.

Some of the references (in both documents) to SQL or GQL draft standards or discussion documents produced in the international standards process link to sites which are only accessible to registered process participants. *Any such document can be made available to any member of LDBC*, but we are not able to make them generally available to the public due to the rules of the international standards process conducted by ISO/IEC JTC1.

Please mail [info@ldbncouncil.org](mailto:info@ldbncouncil.org) to request access or to join LDBC.

This report is published both for general interest, and to provide a reference for other works, where the mathematical definitions and historical explanations may be useful.

The report can be read without any reference to the original discussion paper version. However, it may be useful to know that in producing this report:

- We have changed the title, and removed a prefatory note on standards procedures.
- We introduce new names for paths with no vertices or with only a single vertex, which we feel will help to avoid ambiguities found historically in the research literature relating to the terms "null path" and "empty path" (which we also used previously)
- We also avoid use of the term "null path" because of possible confusion with the SQL concept of a "null value". This latter point arose in a discussion on the original WG013-014 paper, in the INCITS DM32 Ad-Hoc for SQL/PGQ, and in WG3.

#### **Copyright (C) 2021 Linked Data Benchmark Council (LDBC)**

This document is licensed under the terms of the Creative Commons Attribution-NoDerivatives 4.0 International (CC BY-ND 4.0) license, which is linked from here <https://creativecommons.org/licenses/by-nd/4.0/>.

If this document or prior versions thereof have been submitted, or are in the future submitted, to ISO/IEC JTC1 by LDBC for consideration and potential incorporation in the SQL or GQL international standards as part of LDBC's [Category C Liaison with ISO/IEC JTC 1/SC 32 WG3](#), then that submission is additionally covered by the terms of the [DECLARATION FOR PARTICIPANTS IN ISO ACTIVITIES](#).

Any contributor to or recipient of this document who is aware of any patent claim that might be infringed by a specification or standard incorporating some or all of the content of this document is invited to notify the LDBC Steering Committee of the details by sending email to [info@ldbncouncil.org](mailto:info@ldbncouncil.org). At the date of publication LDBC has not been notified of any such patent claims.

**[DOI] material**

## Property graphs and paths in GQL: Mathematical definitions

Alastair Green (LDBC Liaison Representative to WG3)  
Paolo Guagliardo (UK expert, LDBC GQL Formal Semantics WG)  
Leonid Libkin (UK expert, LDBC GQL Formal Semantics WG)

Original **WG3:W13-014** discussion paper, 21 June 2021  
LDBC Technical Report TR-2021-01, **xx** October 2021

### Contents

<b>1. References</b>	<b>4</b>
<b>2. Introduction</b>	<b>5</b>
<b>3. Varying terminology in graph theory</b>	<b>5</b>
<b>4. Graph theoretical and textual definitions of a walk</b>	<b>6</b>
<b>5. An “extracted path”: a partially-oriented walk in a mixed multigraph</b>	<b>6</b>
<b>6. Two equivalent ways of defining a walk mathematically</b>	<b>8</b>
6.1. Sequences and subsequences	8
6.2. Approach I. A sequence alternating between nodes and edges	8
6.3. Approach II. A sequence of nodes, with a corresponding sequence of edges	10
6.4. Translating between Approach I and Approach II	12
<b>7. Proposed concepts section for the SQL/PGQ and GQL specifications</b>	<b>12</b>
<b>8. Proposed informal definitions relating to “path” in SQL/PGQ and GQL</b>	<b>13</b>
<b>9. Appendix: On evolution of terminology</b>	<b>16</b>
<b>10. Appendix: Projection of an extracted path</b>	<b>17</b>
<b>11. Appendix: Diagrams</b>	<b>19</b>
11.1. Walks using Approach I	19
11.2. The same walk using Approach II	20
11.3. Subwalk of the same walk	21
11.4. Trail	22
11.5. Simple path	23

## 1. References

Reference	Title	Author	Date
<a href="#">[9075_8IWD14-16-PGQ]</a>	“Information technology — Database language SQL — Part 16: Property Graph Queries (SQL/PGQ) Working Draft 14	Eds Stefan Plantikow, Stephen Cannan	4 June 2021
<a href="#">[WG3:W10-019r1]</a>	“Aligning the definitions” r1	Stephen Cannan	8 April 2021
<a href="#">[WG3:W10-026r1]</a>	“Definition of a path in SQL/PGQ and GQL”	Alastair Green	21 April 2021
<a href="#">[WG3:W10-027r1]</a>	“Graph Terminology for GQL and SQL/PGQ: Mathematical definitions”	Leonid Libkin	21 April 2021
<a href="#">[WG3:W10-028]</a>	“Input to discussion on path terminology”	Neo4j Query Languages Standards and Research Team	Uploaded 22 April 2021 [misdated 27 November 2020]
<a href="#">[WG3:W11-029]</a>	“<property reference> proposal”	Fred Zemke	19 May 2021
<a href="#">[GraphHomomorphisms]</a>	“Graphs and Homomorphisms”, Oxford Lecture Series in Mathematics and its Applications 28, OUP, pp 1-3, 4, 6]	Pavol Hell and Jaroslav Nešetřil	2004
<a href="#">[Digraphs]</a>	“Digraphs: Theory, Applications and Algorithms”, Springer-Verlag, Berlin.	Jørgen Bang-Jensen and Gregory Gutin	2000 (1st edition), corrected 2007
[Graph Theory]	“Graph Theory”, Addison-Wesley, Reading, Mass.	Frank Harary	1969
[Structural Models]	“Structural models: an introduction to the theory of directed graphs”	Harary, Norman and Cartwright	1965
[Null Graph]	"Is the Null Graph a Pointless Concept?" In <i>Graphs and Combinatorics Conference, George Washington University</i> . New York: Springer-Verlag	Harary, F. and Read, R.	1973

<a href="#">[Empty Graph]</a>	"Empty Graph." From <a href="#">MathWorld</a> --A Wolfram Web Resource.	Weissstein, Eric W.	2002
-------------------------------	---	---------------------	------

## Note on references

Key references for this area in terms of standard concepts and their definitions are [\[GraphHomomorphisms\]](#) and [\[Digraphs\]](#). (The second edition [2009] of [\[Digraphs\]](#) is less useful for this particular area of mixed (partially directed) graphs.)

## 2. Introduction

This paper sets out to

1. Provide a formal mathematical definition of a “property graph” and “path in a property graph”, which we view as a *partially-oriented walk in a mixed multigraph*, a concept which is apparently not fully described in the existing literature of graph theory.
2. Provide definitions of associated concepts such as sub-path, path length, and restricted walks. We allow in these definitions for both 1-vertex (single vertex) walks, with a length of 0, that is consisting only of an isolated vertex), and 0-vertex (zero-vertex) walks, which contain no graph elements, and which can arise in operations that attempt to produce paths.
3. Propose corresponding textual definitions for inclusion in the SQL/PGQ and GQL specifications, and some commentary suitable for a concepts section in either of those documents.

## 3. Varying terminology in graph theory

Many fields have their own distinct cultures. In [Graph Theory] Frank Harary, who played a central role in the development of modern graph theory from the 1950s onwards<sup>1</sup>, observed:

“Most graph theorists use personalized terminology in their books, papers and lectures. In order to avoid quibbling at conferences on graph theory it has been found convenient to adopt the procedure that each man *[sic!]* state in advance the graph theoretic language he would use. Even the very word “graph” is not sacrosanct. Some authors actually define a “graph” as a graph,\* *[Footnote: \*This is most frequently done by the canonical initial sentence, “In this paper we only consider finite undirected graphs without loops or multiple edges.”]* but others intend such alternatives as multigraph, pseudograph, directed graph or network. We believe that uniformity in graphical terminology will never be attained, and is not necessarily desirable.”

---

<sup>1</sup> Harary was the major initial theorist of directed or digraphs, codifying his work in [Structural Models], published in 1965. The term “mixed graph” was first used in 1967 in a paper published by Harary.

As we are developing standards, we are compelled to define a uniform terminology, a goal which Stephen Cannan has pursued in [\[WG3:W10-019r1\]](#).

The tradition of multiple terminologies described by Harary is one obstacle, already familiar from the multiple terms used for nodes, points or vertices, and relationships, lines, arcs and edges.

## 4. Graph theoretical and textual definitions of a walk

For the GQL and SQL standards textual definitions are required. It is desirable that these rest on, or are equivalent to, mathematical definitions that conform as much as possible with standard set theoretical formalizations used in graph theory.

However, a reasonable textual definition may not align exactly with the most acceptable or conventional mathematical definition.

In this case we have an additional factor: although the idea of a mixed multigraph is well-known and studied in graph theory, it appears that the idea of a *partially-oriented path or walk* in such a graph is novel from the standpoint of the published academic literature. Therefore, its mathematical definition will have to extend or synthesize existing partial definitions.

This paper proposes two (equivalent) alternative mathematical approaches, and a set of textual definitions for potential inclusion in the SQL-PGQ (Part 16) working draft (and by editorial extension, to GQL).

We start with a summary of the role of paths in [\[9075\\_8IWD14-16-PGQ\]](#) and how they relate to graph theoretical concepts.

## 5. An “extracted path”: a partially-oriented walk in a mixed multigraph

*Discussion specific to “graph pattern matching” in the draft SQL/PGQ and GQL standards, is highlighted in [bracketed blue text].*

*[Path pattern matching in SQL/PGQ and GQL takes as an input a property graph and a <graph pattern> expression, and produces a set of “reduced” path bindings<sup>2</sup>.*

*A reduced path binding is a sequence of pairs (**element variable** | **anonymous element symbol, graph element**). A graph element is, of course, either a vertex or an edge.*

*A sequence formed by replacing each pair in a reduced path binding with its graph element member alone is a sequence of graph elements, which after removing any duplicate vertices, is called an **extracted path**.*

---

<sup>2</sup> See [“Appendix: Projection of an extracted path”](#) for a detailed summary of this process, which is defined in two key sections of [\[9075\\_8IWD14-16-PGQ\]](#): “9.3 Machinery for graph pattern matching” and “10.2 <graph pattern>”.

An extracted path corresponds to the graph theory concept of a walk  $W$  in a graph  $G$ .]

A walk  $W$  in a graph  $G$  is a second graph, independent of  $G$  but defined with respect to  $G$ , so a vertex in  $W$  or an edge in  $W$  must also be mapped to a vertex or edge in  $G$ . This allows two elements in  $W$  to refer to the same element in  $G$ .

[In our context of an attributed (property) graph in a database, if the identifier, labels, properties and property values of a node or edge in  $W$  are accessed by a client, then they must be those of the element it references in  $G$ , *at the logical point in time when  $G$  is extracted* in a read query (SELECT from GRAPH\_TABLE in SQL/PGQ). Two elements in a walk may have the same element identifier, therefore.]

A walk is a sequence of graph elements with a distinct structure. It can have no members (be a zero-vertex walk/graph); one member (which must be a vertex, making it a single-vertex walk/graph); or an odd number of members 3, 5, 7 ..., in which case the walk starts with a vertex, which is followed by a subsequence of half-edges (each an ordered pair (*edge*, *vertex*)), and therefore ends with a vertex as well, making it a multi-vertex walk.

Property graphs are *multigraphs* (there can be multiple edges between two endpoint vertices) and *pseudographs* (there can be an edge looping from a vertex to itself). They are *mixed* (partially directed) *graphs*, as an edge can be *undirected*, or can have tail and head vertices, in which case it is *directed* from tail to head (or equivalently is *dominated* by the head vertex).

A walk  $W$  in a mixed graph  $M$  traverses vertices, (undirected) edges and arcs (directed edges). Such a walk therefore can be defined to contain information about the role that a vertex plays as an endpoint: whether it is neutral with respect to the direction of an edge, or acts as the head or tail of an arc (as directed edges are sometimes termed). This information gives us the *orientation* of an edge *with respect to the walk*, which (as a sequence) is ordered, and therefore has a first and last element. (A vertex can therefore act as the endpoint of zero, one or two edges in a walk).

An edge in  $W$  may be *unoriented* (if that edge is undirected in  $M$ ), or it may be *oriented forwards* if the edge is directed and its tail endpoint occurs before its head endpoint in the walk, or *backwards* if the reverse is true.

[The consumer of an extracted path should be able to find the identity and content of each graph element in the path viewed as a walk, its position in the path, and its orientation (if an edge) or role in the orientation of each edge for which it is an endpoint (if a vertex).

An extracted path in SQL/PGQ or GQL is therefore an example of] a walk  $W$  in  $G$ , that is, a **partially-oriented walk in a property graph**, which is an attributed, mixed multi/pseudograph<sup>3</sup>.

---

<sup>3</sup> In the academic literature we have found work on wholly-oriented walks in mixed multigraphs (described and referenced in [Digraphs], 8.9 “Orientations of Mixed Graphs”), but as the authors explain, for their purposes “when we speak of orienting a mixed (multi)graph this means that we assign an orientation to every edge and leave the original arcs unchanged”. In our case we wish to leave the edges unoriented, and our orientation of the graph (walk) is therefore partial.

A walk where each edge in the path is also a distinct edge in the graph is called a trail; a walk where each vertex in the path is also a distinct vertex in the graph is called a simple path.

## 6. Two equivalent ways of defining a walk mathematically

There are (at least) two valid approaches to defining a walk in mathematical terms.

To start, we define our notation for sequences and subsequences.

### 6.1. Sequences and subsequences

A *sequence*  $s$  of length  $n \in \mathbb{N}_0$  over a set  $\mathcal{A}$  is a function  $s : \{i \in \mathbb{N}_1 \mid 1 \leq i \leq n\} \rightarrow \mathcal{A}$ , often written by listing its elements explicitly, e.g.,  $(s_1, \dots, s_n)$ . Given  $s$ , by  $s[i]$  we denote its  $i$ -th element (in the previous example,  $s_i$ ). If  $n = 0$ , this is the *empty* sequence.

Given a sequence  $s$  of length  $n$ , and bounds  $b$  and  $e$  (beginning and end) with  $1 \leq b \leq e \leq n$ , the *subsequence*  $s'[b\dots e]$  of  $s$  is the sequence  $s' = (s[b], \dots, s[e])$ , that is,  $s'[1] = s[b]$ ,  $s'[2] = s[b+1]$ ,  $\dots$ ,  $s'[e-b+1] = s[e]$ . The empty sequence is a subsequence of every sequence.

For example, in a sequence  $ab = (a_1, b_1, a_2, b_2, a_3)$  we have  $ab[1] = a_1$ ,  $ab[2] = b_1$ ,  $ab[3] = a_2$ ,  $ab[4] = b_2$ ,  $ab[5] = a_3$ . The subsequence  $ab' = ab[2\dots 4]$  is  $(b_1, a_2, b_2)$ . So, in this sequence,  $ab'[1] = ab[2] = b_1$ ,  $ab'[2] = s[3] = a_2$  and  $ab'[3] = s[4] = b_2$ .

### 6.2. Approach I. A sequence alternating between nodes and edges

Approach I rests on the concept of a sequence of graph elements alternating between nodes and edges. Cannan, [\[WG3:W10-019r1\]](#), Libkin, [\[WG3:W10-027r1\]](#), Neo4j [\[WG3:W10-028\]](#). This approach is used in e.g. [\[Digraphs\]](#) by Bang-Jensen and Gutin.

In [\[9075\\_8|WD14-16-PGQ\]](#) this is formulated textually in “9.3 Machinery for graph pattern matching”, p28 as:

“... either a zero-length string or it is a path of  $PG$  [the pure property graph]; it begins with a vertex, it alternates between vertices and edges, each edge connects the vertex before and after it in the extracted path, and it ends with a vertex.”

(although this treats an empty sequence [one with no elements at all] as not being a path).

[\[WG3:W10-027r1\]](#) rendered this approach formally, and we have expanded the definitions given there, as follows.

A property graph is a tuple  $G = \langle N, E_d, E_u, \lambda, \text{ENDPOINTS}, \text{SRC}, \text{TGT}, \text{DV} \rangle$  where

- $N$  is a set of node ids used in  $G$ , coming from a set  $\mathcal{N}$  of node ids;
- $E_d$  is a set of directed edge ids used in  $G$ , coming from a set  $\mathcal{E}_d$  of directed edge ids;
- $E_u$  is a set of undirected edge ids used in  $G$ , coming from a set  $\mathcal{E}_u$  of undirected edge ids;



- $\lambda$  is a labeling function;  $\lambda : N \cup E_d \cup E_u \rightarrow 2^{\mathcal{L}}$  associates with each of the ids a set (possibly empty) of labels from the set  $\mathcal{L}$  of labels<sup>4</sup>;
- **ENDPOINTS**:  $E_u \rightarrow 2^N$  associates with each edge is  $e \in E_u$  a set of its endpoints such that the cardinality of **endpoints**( $e$ ) is 1 or 2;
- **SRC**, **TGT**:  $E_d \rightarrow N$  associate with each directed edge id its source and target;
- **DV**:  $N \cup E_d \cup E_u \rightarrow V$  associates with each id of a node or an edge some data value from a set  $V$ ; normally this would be a set of key-value pairs but for definitions of paths it is not important what this set is;
- all sets  $N$ ,  $E_d$ ,  $E_u$ ,  $\mathcal{L}$ ,  $V$  are disjoint.

Notes:

A graph itself can have data values associated with it. This can be modelled by adding an element of  $V$  to the tuple defining  $G$ . For the purposes of defining paths this makes no difference.

Other possible graph theory terminology for undirected edges: just *edge*, or *line*.

Other possible graph theory terminology for directed edges: *arc*, or *link*.

Two nodes  $v$  and  $v'$  are *adjacent* if either **endpoints**( $e$ ) =  $\{v, v'\}$  for some  $e \in E_u$ , or  $\{v, v'\} = \{\text{SRC}(e), \text{TGT}(e)\}$  for some  $e \in E_d$ . Spelling this out, to help with plain language definitions, means either  $v = \text{SRC}(e)$ ,  $v' = \text{TGT}(e)$  or  $v' = \text{SRC}(e)$ ,  $v = \text{TGT}(e)$  for some  $e \in E_d$ .

**Definition I.1 (Path)** A path in  $G$  is a sequence  $\pi = (v_1, e_1, v_2, e_2, \dots, e_{n-1}, v_n)$  where:

1. each  $v_i$  is a node id from  $N$  and each  $e_j$  is an edge id from  $E_d \cup E_u$ ;
2.  $n \geq 0$ ;
3. for every  $j$  between 1 and  $n-1$ , one of the following is true:
  - $e_j \in E_u$  and **ENDPOINTS**( $e_j$ ) =  $\{v_j, v_{j+1}\}$  (an *undirected edge*);
  - $e_j \in E_d$  and  $\text{SRC}(e_j) = v_j$  and  $\text{TGT}(e_j) = v_{j+1}$  (a *forward edge* or *forward arc*);
  - $e_j \in E_d$  and  $\text{TGT}(e_j) = v_j$  and  $\text{SRC}(e_j) = v_{j+1}$  (a *backward edge* or *backward arc*).

Note: This definition allows for a walk consisting of an isolated, single vertex (an *empty path*) or a *null path* (one with no vertices), by relying on the fact that if the number of vertices  $n = 0 \mid 1$ , then there is no positive integer value  $j$  between 1 and  $n-1 = 0 \mid -1$ , and therefore there is no member of the set of edges in the sequence.

We now define the length of a path, a sub-path, and associated concepts:

**Definition I.2 (Length of a path)**

The length of a path is  $n-1$ .

---

<sup>4</sup> Various adjustments are possible, e.g. by saying  $|\lambda(e)| = 1$  for  $e \in E_d \cup E_u$  to force single edge labels etc.

**Definition I.3 (Zero-edge or single-vertex path:  $n=1$ )**

A path is zero-edge, or equivalently single-vertex, if  $n=1$ ; such a path has a length of 0.

**Definition I.4 (Zero-vertex path:  $n=0$ )**

A path is zero-vertex if  $n=0$ ; such a path has a length of  $-1$ .

**Definition I.5 (Sub-path)**

Given a path  $\pi = (v_1, e_1, v_2, e_2, \dots, e_{n-1}, v_n)$ , its subpaths are all subsequences of  $\pi$  which are themselves paths, i.e., the empty sequence or  $(v_i, e_i, \dots, e_{j-1}, v_j)$  whenever  $1 \leq i \leq j \leq n$ .

To spell this out,  $(v_i, e_i, \dots, e_{j-1}, v_j) = (\pi[2i-1], \pi[2i], \pi[2i+1], \dots, \pi[2j-1])$ .

Thus, a single node (empty path) is a subpath (when  $i = j$ ) and the null path is always a subpath corresponding to the empty subsequence.

**Definitions I.6 (Special types of paths)**

A path  $\pi = (v_1, e_1, v_2, e_2, \dots, e_{n-1}, v_n)$  is **simple** if no node repeats except perhaps the first and the last one. That is, if  $\{i, j\}$  is not equal to  $\{1, n\}$  and  $i$  and  $j$  are different, then  $v_i$  and  $v_j$  are different.

A path is **acyclic** if it is simple and  $v_1$  and  $v_n$  are different.

A path is a **trail** if  $e_i$  and  $e_j$  are different whenever  $i$  and  $j$  are different.

A path is **closed** if  $v_1 = v_n$ .

### 6.3. Approach II. A sequence of nodes, with a corresponding sequence of edges

The second approach defines a walk as a sequence of 0, 1 or more vertices, where each consecutive pair in a sequence of 2 or more vertices maps to an edge in a corresponding sequence of 1 or more edges. The sequence of edges corresponding to a vertex sequence of length one, or nought, is empty.

The interface for Neo4j’s client driver libraries for [paths](#) shows this approach:

<code>Iterable&lt;Node&gt;</code>	<code>nodes()</code> Create an iterable over the nodes in this path, nodes will appear in the same order as they appear in the path.
<code>Iterable&lt;Relationship&gt;</code>	<code>relationships()</code> Create an iterable over the relationships in this path.

First we define a graph using the conventional terminology often seen when treating directed graphs, where undirected lines are edges, and directed ones are arcs. This definition also removes graph element identifiers, adopting an alternative approach of defining multisets of edges and arcs (taken together, a multiset of relationships).

A **property graph** is a partially directed mixed pseudo-multigraph  $M = (V, E, A)$ , where

$V$  is a set of vertices;

$E$  is a multiset of edges where each  $e \in E$  is an unordered pair of endpoints  $\{u, v\} \subseteq V$ ;

$\mathcal{A}$  is a multiset of arcs,  $V \times V$ . Each  $a \in \mathcal{A} = (u, v)$  is an ordered pair of endpoints, where the arc is directed from  $u$  (its tail) to  $v$  (its head).

We term  $E \cup \mathcal{A}$  as the multiset of relationships  $\mathcal{R}$ : for any pair of endpoints  $\{u, v\}$  there may be many parallel edges and parallel arcs (which may be of opposing directions);

$\mathcal{M}$  and each member of  $V \cup \mathcal{R}$  is mapped to a set of attributes (a disjoint union of a set of labels and properties).

A **partially-oriented walk**  $W(n, m)$  in  $\mathcal{M}$  is a pair of sequences  $\langle v, r \rangle$  where:

$n \geq 0$  and  $m = \max(n-1, 0)$ ;

$v = (v_i)_{1 \leq i \leq n}$  is a sequence of  $n$  vertices from  $V$ ;

$r = (r_i)_{1 \leq i \leq m}$  is a sequence of  $m$  relationships from  $\mathcal{R}$  such that each  $r_i$  is one of the following:

1. an edge  $e \in E$  whose endpoints are  $\{v_i, v_{i+1}\}$
2. a forward arc  $f \in \mathcal{A}$  whose tail is  $v_{i+1}$ , and whose head is  $v_i$
3. a backward arc  $b \in \mathcal{A}$  whose tail is  $v_i$ , and whose head is  $v_{i+1}$

Note: as a walk is a graph (separate from the graph  $\mathcal{M}$  from which it is derived) we can use the conventional notation  $W(n, m)$  to express the order (number of vertices) and size (number of edges) in this graph. This allows us to easily note the difference between a 0-vertex walk and a 1-vertex (equivalently, a 0-edge) walk, both of which have a size of 0.

This approach is suggested by the view of a path as centring on a sequence of vertices, used in e.g. [\[GraphHomomorphisms\]](#) by Hell and Nešetřil, and in Green [\[WG3:W10-026r1\]](#).

We now introduce consequential definitions:

### Definition II.2 (Length and size of a walk)

The length of a walk  $W(n, m)$  is  $n-1$ , while its size is the number of its edges  $m$ .<sup>5</sup>

### Definition II.3 (Single-vertex, or zero-edge, walk)

If  $n = 1$  then  $m = 0$ : a walk  $W(1, 0)$  is an single-vertex walk (with length 0), containing only an isolated vertex.

### Definition II.4 (Zero-vertex walk)

If  $n = 0$  then  $m = 0$ : a walk  $W_\emptyset = W(0, 0)$  is a zero-vertex walk (with length  $-1$ ).

### Definition II.5 (Subwalk)

A subwalk of a walk  $W(n, m) = \langle v, r \rangle$  is either a zero-vertex walk or the walk  $\langle v', r' \rangle$  where, for bounds  $b$  and  $e$  with  $1 \leq b \leq e \leq n$ , the sequence  $v'$  is  $v[b...e]$  and the sequence  $r'$  is  $r[b...e-1]$ .

<sup>5</sup> Graphs conventionally have sizes, and paths have lengths (both being the number of edges). In our context, where we may have no edges and no vertices in a zero-vertex (sub-)path it is useful to distinguish the size of a path *qua* graph, from its length *qua* path.

## Definitions II.6 (Special types of walks)

Recall that sequences  $\mathbf{v}$  and  $\mathbf{r}$  are functions  $\mathbf{v} : \{1, \dots, n\} \rightarrow V$  and  $\mathbf{r} : \{1, \dots, m\} \rightarrow R$ .

A walk is **acyclic** if  $\mathbf{v}$  is injective.

A walk is **simple** if the restrictions of  $\mathbf{v}$  to both of  $\{1, \dots, n-1\}$  and  $\{2, \dots, n\}$  are injective functions.

A walk is **closed** if  $\mathbf{v}(1) = \mathbf{v}(n)$  (which is equivalent to  $\mathbf{v}[1] = \mathbf{v}[n]$ ).

A walk is a **trail** if the function  $\mathbf{r}$  is injective.

The diagrams in [“Appendix: Diagrams”](#) illustrate a walk, a subwalk and these special walks.

## 6.4. Translating between Approach I and Approach II

Both of these formal definitional approaches describe the same class of mathematical object.

If we have a path

$$\pi = (\mathbf{v}_1, \mathbf{e}_1, \mathbf{v}_2, \mathbf{e}_2, \dots, \mathbf{e}_{n-1}, \mathbf{v}_n)$$

then the equivalent partially-oriented walk is

$$\begin{aligned} \mathbf{W}(n, m)_\pi &= \langle \mathbf{v}, \mathbf{r} \rangle \text{ where} \\ m &= \max(n-1, 0) \\ \mathbf{v} &= (\pi[1], \pi[3], \dots, \pi[2n-1]) \\ \mathbf{r} &= (\pi[2], \pi[4], \dots, \pi[2n]). \end{aligned}$$

Given a walk

$$\mathbf{W}(n, m)$$

the corresponding path

$$\pi_{\mathbf{W}(n, m)} = (\mathbf{v}[1], \mathbf{r}[1], \mathbf{v}[2], \mathbf{r}[2], \dots, \mathbf{r}[n-1], \mathbf{v}[n]).$$

The following sections propose textual definitions and descriptions drawing on both of these approaches. We start by proposing some wording on concepts.

## 7. Proposed concepts section for the SQL/PGQ and GQL specifications

The following text could be included in a concepts section of both the SQL/PGQ and GQL specifications.

“A mathematical definition of a *mixed (partially directed)* graph is given in “Digraphs: Theory, Applications and Algorithms”, Jørgen Bang-Jensen and Gregory Gutin. 2007, Springer-Verlag, Berlin. Section “8.9 Orientations of Mixed Graphs”. In summary:

Let  $\mathbf{M} = (\mathbf{V}, \mathbf{A}, \mathbf{E})$  be a mixed graph. Let  $\mathbf{G} = (\mathbf{V}, \mathbf{E})$  be the undirected part and let  $\mathbf{D} = (\mathbf{V}, \mathbf{A})$  be the directed part of  $\mathbf{M}$ .

The sets of vertex pairs that define the (undirected) edges in  $\mathbf{E}$  and the (directed) arcs in  $\mathbf{A}$  can overlap, so at most there can be an edge and two arcs between any two vertices in  $\mathbf{V}$ .

The authors add:

A mixed multigraph is the same as a mixed graph, except that we allow parallel arcs and parallel edges as well as arcs that are parallel to edges.

A mathematical definition of an oriented walk in a *directed* graph is given in the “Introduction” to Hell, Pavol, and Jaroslav, Nešetřil. 2004. [Graphs and Homomorphisms](#). OUP Oxford, pp 4, 6.<sup>6</sup>

... an oriented walk in a digraph  $G$  is a sequence of vertices  $v_0, v_1, \dots, v_k$  of  $G$  such that  $v_{i-1}$  and  $v_i$  are adjacent in  $G$ , for each  $i = 1, 2, \dots, k$ . An arc  $v_{i-1} v_i \in E(G)$  is called a forward arc of the walk, and an arc  $v_i v_{i-1}$  is called a backward arc of the walk.

In a database property graph a walk in a *mixed* (partially directed) *multigraph* may also a) traverse parallel arcs and edges, and b) traverse an edge which is undirected: the walk is therefore only *partially oriented*.

In this standard the unqualified term “path” is used as a convenient and familiar synonym for a ***partially-oriented walk in a mixed multigraph*** which manifests itself, for example, as an extracted path resulting from a graph pattern match.

The fact that an edge may be a self-loop (that the graph is a pseudograph) does not affect the concept of a path or walk.”

## 8. Proposed informal definitions relating to “path” in SQL/PGQ and GQL

*Definitions 3, 4, 5 and the second sentence of 6 paraphrase or relate to terms specific to “graph pattern matching” in the draft SQL/PGQ and GQL standards, and are highlighted in [bracketed blue text].*

### 1. Property graph

A graph which is

- a. *attributed*: the graph, and each of its vertices and edges, has a (possibly empty) set of attributes, being the disjoint union of a set of labels and a set of properties
- b. *partially directed* (mixed), so each edge may be directed from a source vertex to a destination vertex, or may be undirected (neither of its endpoints is a source or a destination)
- c. a *multigraph*, allowing multiple edges between two vertices
- d. a *pseudograph*, allowing loops, which are edges which join a vertex to itself.

### 2. Partially-oriented walk

A walk (as commonly understood in graph theory) in a property graph, whose edges are either unoriented—if an edge has a corresponding edge in the graph which is undirected—or oriented, when an edge reflects the direction of a corresponding directed edge in the graph. The orientation of an oriented edge is defined with respect to the order of the walk, which is a sequence from an initial vertex to a final vertex, if it is not an empty sequence.

---

<sup>6</sup> <https://paperpile.com/shared/TqO0j6>

See also **Path**, **<Path>Unoriented edge**, **<Path>Oriented Edge**, **<Path>Null path**.

3. **[Path binding]**

A possibly empty sequence resulting from matching a graph pattern in a property graph, whose terms are each a pair, of the following kinds:

1. A bracket pair, whose first and second members are both either an opening or a closing bracket symbol
2. A subpath delimiter pair, whose first and second members are both the same start or end subpath symbol, unique in the path binding, allowing a subpath to be delimited
3. A vertex pair, whose first member is a vertex variable or an anonymous vertex symbol and whose second member is a vertex in the property graph
4. An edge pair, whose first member is an edge variable or an anonymous edge symbol,, and whose second member is an edge in the property graph

and whose first and last terms must both be vertex pairs.]

4. **[<Path binding>Annotated path]**

A sequence of the second members of the pairs in a binding path, in the same order as the binding path sequence.]

5. **[<Annotated path>Extracted path]**

An extracted path is a subsequence of an annotated path produced by working from beginning to end, removing bracket symbols, subpath symbols and consecutive (duplicate) vertices. Each remaining edge is preceded by a vertex and followed by a vertex, which are its endpoints in the property graph from which the path was extracted. An extracted path may contain only one vertex or no vertices, in which cases it contains no edges. The edges of an extracted path are either oriented or unoriented.]

6. **Path**

A path is a partially-oriented walk in a property graph.[ A path is a sequence with the structure and characteristics of an extracted path.]<sup>7</sup>

7. **<Path>Associated property graph**

The graph with respect to which a walk is defined is the associated property graph of that walk.

8. **<Path>Oriented edge**

---

<sup>7</sup> [It is possible that a path may be created by other means than pattern matching, for example in the course of processing a GQL INSERT or DELETE statement. A statement making this explicit might be appropriate as a note in the GQL specification.]

An edge in a path, which is a directed edge in its associated property graph, is *forward-oriented* if the vertex to its left in the path is its source endpoint, or *backward-oriented* if the vertex to its left is its destination endpoint.

9. **<Path>Unoriented edge**

An *unoriented* edge in a path is an undirected edge in its associated property graph.

10. **<Path>Length**

The length of a path is one less than the number of its vertices. If a path contains edges then its length is the number of its edges.

11. **<Path>Zero-vertex path**

A path with no vertices (and whose length is therefore  $-1$ ).

12. **<Path>One-vertex path**

A path with one vertex (and whose length is therefore 0).

13. **<Path>Closed**

A path is closed (synonymously, is a cycle) if its first and last members are the same vertex in its associated property graph.

14. **<Path>Subpath**

A subpath is a path which is a contiguous subsequence of another path, and may be an empty sequence, i.e a zero-vertex path.

15. **<Path>Walk**

The unqualified term walk is a synonym for path.

16. **<Path>Trail**

A walk such that no two of its edges are the same edge in its associated property graph.

17. **<Path>Simple path**

A walk such that only the first and last terms may be the same vertex in its associated property graph.

18. **<Path>Acyclic**

A path is acyclic if no two of its vertices are the same vertex in its associated property graph. It is therefore a simple path.

## 9. Appendix: On evolution of terminology

The original version of this paper was discussed in meetings of the U.S. standards committee INCITS DM32 Ad-Hoc for SQL/PGQ (Property Graph Query), and of the international standards committee ISO/IEC JTC1 SC32/WG3, in July 2021. It was pointed out that the term “null” has a very strong connotation in SQL of a null value (of a cell in a table or of an expression), and that the term “null walk” could be confusing.

On reflection, we can also see that walks with no content, and walks with one vertex but no edges could both be described as “empty”, if by that we mean edgeless. To make things worse, there is a standing ambiguity in the academic literature on the meaning of “null graph”, which is reflected in this entry in the Wolfram MathWorld online dictionary of mathematics: [\[Empty Graph\]](#).

There is an old debate, summarized in [Null Graph] about the validity of the concept of a graph (of which a walk is a special case) with no vertices.

In the context of path pattern matching in a graph, a pattern may yield no matches (producing a graph with no elements). If a path pattern contains a sub-path pattern, and the sub-path pattern has no matches, then the path pattern as a whole may still produce matches. In this context we do need to define a matching operation that can yield a walk that is completely void (has no elements). For this reason we see value in defining special walks with no vertices and with one vertex.

Therefore we have decided to avoid the terms “empty walk” and “null walk” altogether, in favour of the terms “zero-vertex walk” and “single-vertex walk”, which are unambiguous and self-descriptive.

Viewed as amendments to our original definitions (using Approach II, Section x.3), where we define

**A partially-oriented walk  $W(n, m)$  in  $M$  ...**

we see the following changes to supplementary definitions:

**Definition (1) (~~Empty walk~~ Single-vertex or zero-edge walk)**

If  $n = 1$  then  $m = 0$ : a walk  $W(1, 0)$  is a single-vertex ~~empty~~ walk (with length 0), containing only an isolated vertex.

**Definition (2) (~~Null walk~~ Zero-vertex walk)**

If  $n = 0$  then  $m = 0$ : a walk  $W_\emptyset = W(0, 0)$  is a zero-vertex ~~null~~ walk (with length  $-1$ ).

We could then further say:

**Definition (3) (Multi-vertex walk)**

If  $n > 1$  then  $m = n - 1$ , and a walk  $W(n, m)$  is a multi-vertex walk (with length  $m$ ).

We can equivalently talk about 1-vertex, 0-vertex, 2-vertex (edge) and M-vertex walks.



## 10. Appendix: Projection of an extracted path

*This appendix describes the process of “graph pattern matching” in the draft SQL/PGQ and GQL standards. The concepts and definitions in this report are stated independently of this motivating use-case of “path extraction”.*

An extracted path is produced (logically) by the application of a sequence of procedures operating on a **pure property graph** and a <path pattern list> from a **<graph pattern>**.

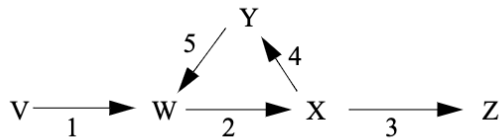
These steps are described in two key sections of [\[9075\\_8IWD14-16-PGQ\]](#): “9.3 Machinery for graph pattern matching” and “10.2 <graph pattern>”

On p. 36 of “<property reference> proposal” [\[WG3:W11-029\]](#), Fred Zemke gives a useful example which helps to illustrate the projection of an extracted path. It is repeated here, with emphasis added in a couple of places:

Consider the pattern

$$[(A) \rightarrow [(B) - [E] \rightarrow (C) \text{ WHERE } E.X = A.X + D.X] + \rightarrow (D)]\{2\}$$

to be evaluated in this graph:



I will use  $[_1]_1$  as the bracket symbols for the outer <parenthesized path pattern expression>, and  $[_2]_2$  for the inner group.

The candidate solution I consider is  $V \rightarrow W \rightarrow X \rightarrow Y \rightarrow W \rightarrow X \rightarrow Z$ , with the following **path binding**:

$$[_1 A - [_2 B E C]_2 - D]_1 \quad [_1 A - [_2 B E C]_2 - D]_1 \\
[_1 V 1 [_2 W 2 X]_2 4 Y]_1 \quad [_1 Y 5 [_2 W 2 X]_2 3 Z]_1$$

The first sequence shown is the **word** of the path binding; the second is its **annotated path**.

A path binding arises when an implementation finds paths in the pure property graph that match a path pattern. There is a list of path patterns in a graph pattern: each path pattern has a corresponding path binding.

This structure binds element variables in the pattern to graph elements; it also introduces anonymous element symbols associated with graph elements that are not bound by an element variable, and indexed brackets to denote the nesting level and inclusion of bracketed expressions.

A path pattern can contain a sub-path variable, which generates a matched pair of subpath start and end delimiters with a unique identity corresponding to the variable, which surround a subword of the word, and

a corresponding contiguous subsequence of the annotated path. Subpaths are not illustrated in the example.

A path binding can also be defined, therefore as a sequence of pairs: in the example:

$$(([_1, _1], (A, v), (-, 1), \dots, (]_1, ]_1))$$

Brackets are paired with brackets; element variables like  $A$  are paired with a vertex or an edge, and anonymous element symbols  $()$  and  $-$  are paired with vertices and edges respectively. Subpath start and end identifiers are paired with themselves, in the manner of brackets.

An annotated path is the second sequence shown in the example,

$$[_1 v 1 \quad [_2 w 2 x ]_2 \quad 4 y ]_1 \quad [_1 y 5 \quad [_2 w 2 x ]_2 \quad 3 z ]_1$$

that is, a sequence of the second element of each pair in the sequence previously shown, in the same order:

$$([_1, v, 1, \dots, ]_1)$$

An **extracted path** is an annotated path, with all brackets and sub-path delimiters removed:

*After step 1:*  $v \ 1 \ w \ 2 \ x \ 4 \ y \ y \ 5 \ w \ 2 \ x \ 3 \ z$

and all consecutive vertices replaced by a single vertex.

*After step 2:*  $v \ 1 \ w \ 2 \ x \ 4 \ y \ 5 \ w \ 2 \ x \ 3 \ z$

An extracted path conforms to the rule 9.3 (14), that is, it is

... either a zero-length string or it is a path of  $PG$  [the pure property graph]; that is, it begins with a vertex, it alternates between vertices and edges, each edge connects the vertex before and after it in the extracted path, and it ends with a vertex.

The function *REDUCE* defined in 9.3 removes brackets and sub-path delimiters from a path binding.

A “reduced path binding” is therefore a sequence of pairs, and a sequence of the second elements in those pairs appears to be equivalent to an annotated path, after the application of step 1 above, or to an extracted path, before the application of step 2. It may be that closer reading of [\[9075\\_8|WD14-16-PGQ\]](#) will show that this apparent anomaly is our misunderstanding. This has no bearing on the definition of an extracted path, which can be obtained by applying steps 1 and 2 to an annotated path, or step 2 to a sequence of the second elements of the pairs of a reduced path binding.

## 11. Appendix: Diagrams

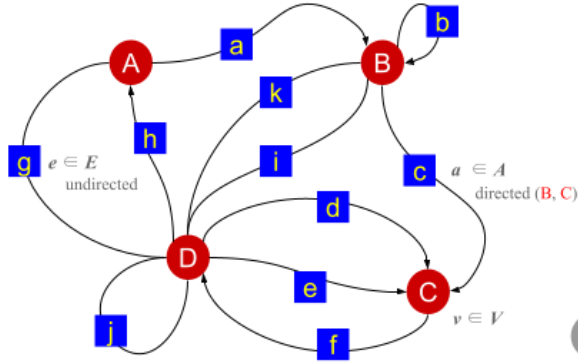
### 11.1. Walks using Approach I

This diagram uses the definition of a mixed multigraph for Approach II, to illustrate the definition of a path using Approach I.

A mixed multigraph  $M = \langle V, E, A \rangle$ ,  $R = E \cup A$

$$V(M) = \{A, B, C, D\}$$

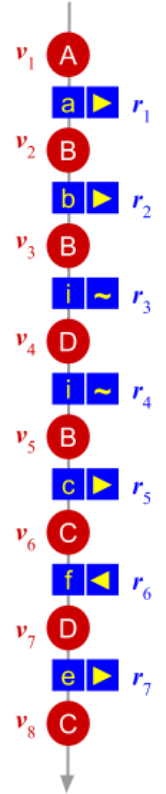
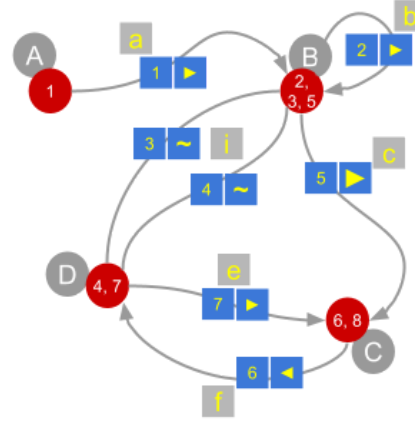
$$R(M) = \{a, b, c, d, e, f, g, h, i, j, k\}$$



A path is a sequence of alternating vertices and edges from  $M$

$$\pi = (v_1, r_1, v_2, r_2, \dots, r_{n-1}, v_n)$$

whose edges may be oriented, if they are directed in the underlying graph.



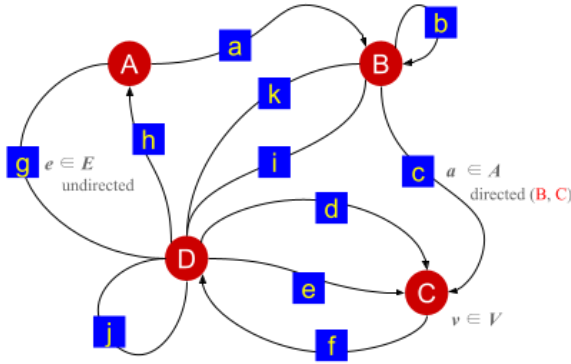
## 11.2. The same walk using Approach II

This diagram shows the equivalent definition of a walk in  $M$  using Approach II.

A mixed multigraph  $M = \langle V, E, A \rangle, R = E \cup A$

$$V(M) = \{A, B, C, D\}$$

$$R(M) = \{a, b, c, d, e, f, g, h, i, j, k\}$$



$$W(8, 7) \text{ in } M = \langle v = (v_n)_{n=8}, r = (r_m)_{m=n-1} \rangle$$

$$v = (v_1, v_2, v_3, v_4, v_5, v_6, v_7, v_8) = (A, B, B, D, B, C, D, C)$$

$$r = (r_1, r_2, r_3, r_4, r_5, r_6, r_7) = (a, b, i, i, c, f, e)$$

There may be no relationships in a walk:  $m$  does not always equal  $n - 1$ .

$$m = \max(n - 1, 0).$$

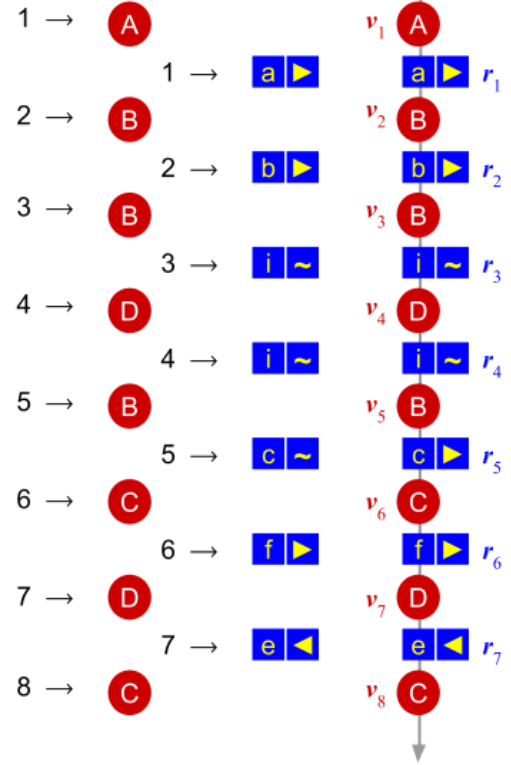
$W(0, 0) = W_\emptyset$  is a *zero-vertex walk*.

$W(1, 0)$  is a *one-vertex walk*.

A walk  $W(n, m)$  in  $M = \langle (v_i)_{0 \leq i \leq n}, (r_j)_{0 \leq j \leq m} \rangle$

$$(v_i)_{0 \leq i \leq n} : \mathbb{N} \rightarrow V(M)$$

$$(r_j)_{0 \leq j \leq m} : \mathbb{N} \rightarrow R(M)$$



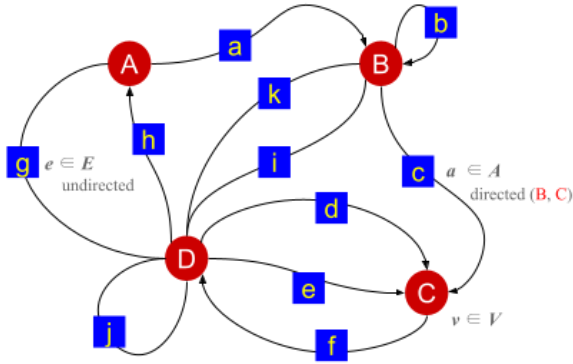
### 11.3. Subwalk of the same walk

This diagram shows a subwalk (which is a walk) of  $W(n, m)$  in  $M$ .

A mixed multigraph  $M = \langle V, E, A \rangle$ ,  $R = E \cup A$

$V(M) = \{A, B, C, D\}$

$R(M) = \{a, b, c, d, e, f, g, h, i, j, k\}$



For a walk  $w = W(8, 7)$  in  $M$

$w = \langle v = (v_n)_{n=8}, r = (r_m)_{m=n-1} \rangle$

$v = (v_1, v_2, v_3, v_4, v_5, v_6, v_7, v_8) = (A, B, B, D, B, C, D, C)$

$r = (r_1, r_2, r_3, r_4, r_5, r_6, r_7) = (a, b, i, i, c, f, e)$

the subwalk

$sw_w = w[2...4] = \langle (v' = v[2...4], r' = r[2...3]) \rangle$

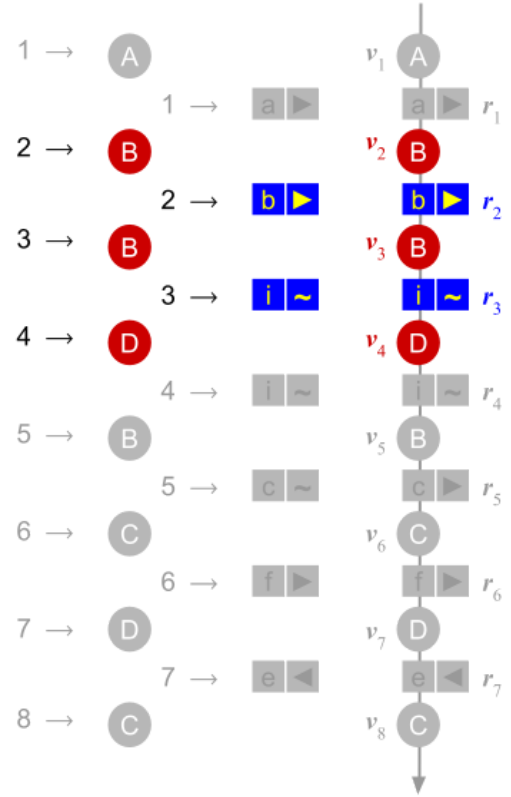
$v' = (v'[1...3]) = (v[2...4]) = (v_2, v_3, v_4) = (B, B, D)$

$r' = (r'[1...2]) = (r[2...3]) = (r_2, r_3) = (b, i)$

$W_\emptyset$  is a subwalk of every walk.

A bounded subwalk  $W(n, m)[b...e]$  in  $M$

$= \langle (v'_i)_{0 \leq i \leq n} [b...e], (r'_j)_{0 \leq j \leq m} [b...e-1] \rangle$ .



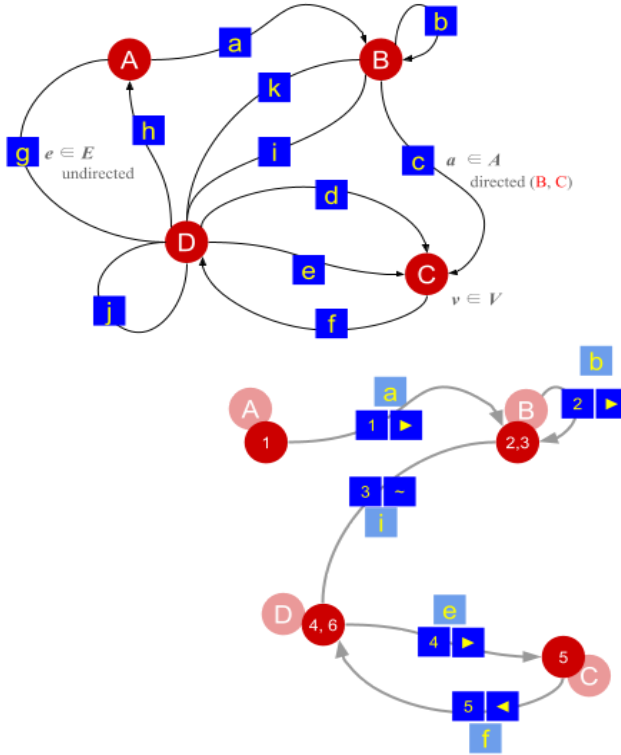
## 11.4. Trail

This diagram shows a trail in  $M$ .

A mixed multigraph  $M = \langle V, E, A \rangle$ ,  $R = E \cup A$

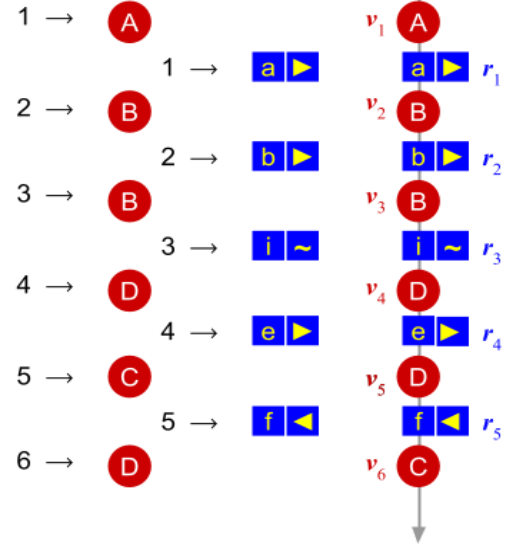
$V(M) = \{A, B, C, D\}$

$R(M) = \{a, b, c, d, e, f, g, h, i, j, k\}$



A trail is a walk  $W(n, m)$  with no repeated relationships: the function that defines the sequence of relationships is injective (one-to-one).

$(r_j)_{0 \leq j \leq m} : \mathbb{N} \hookrightarrow R(M)$



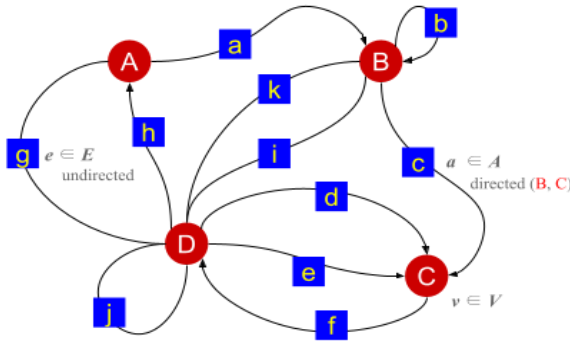
## 11.5. Simple path

This diagram shows a simple path in  $M$ .

A mixed multigraph  $M = \langle V, E, A \rangle$ ,  $R = E \cup A$

$V(M) = \{A, B, C, D\}$

$R(M) = \{a, b, c, d, e, f, g, h, i, j, k\}$



A simple path is a walk  $W(n, m)$  with no repeated vertices (and therefore, no repeated edges): the function that defines the sequence of vertices is injective (one-to-one).

$(v_i)_{0 \leq i \leq n} : \mathbb{N} \hookrightarrow V(M)$

1  $\rightarrow$  A  
2  $\rightarrow$  B  
3  $\rightarrow$  D  
4  $\rightarrow$  C

1  $\rightarrow$  a  $\rightarrow$   
2  $\rightarrow$  i  $\sim$   
3  $\rightarrow$  e  $\rightarrow$

