



HAL
open science

Algorithmic and complexity aspects of path computation in multi-layer networks

Mohamed Lamine Lamali, Nasreddine Fergani, Johanne Cohen

► **To cite this version:**

Mohamed Lamine Lamali, Nasreddine Fergani, Johanne Cohen. Algorithmic and complexity aspects of path computation in multi-layer networks. 2018. hal-01706362v1

HAL Id: hal-01706362

<https://hal.science/hal-01706362v1>

Preprint submitted on 11 Feb 2018 (v1), last revised 1 Nov 2018 (v2)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Algorithmic and complexity aspects of path computation in multi-layer networks

Mohamed Lamine Lamali¹, Nasreddine Fergani², Johanne Cohen³,

¹Université de Bordeaux – LaBRI. France.

²MathWorks. UK.

³LRI, Univ. Paris-Sud, CNRS, Université Paris-Saclay. France.

mohamed_lamine.lamali@u-bordeaux.fr johanne.cohen@lri.fr

Carrier-grade networks comprise several layers where different protocols coexist. Nowadays, most of these networks have different control planes to manage routing on different layers, leading to a suboptimal use of the network resources and to additional operational costs. However, some routers are able to encapsulate, decapsulate, and convert protocols, and act as a liaison between these layers. A unified control plane would be useful to optimize the use of the network resources and to automate the routing configurations. Software-Defined Networking based architectures, offer an opportunity to design such a control plane. One of the most important problems to deal with in this design is the path computation process. Classical path computation algorithms cannot resolve the problem as they do not take into account encapsulations and conversions of protocols. In this paper, we propose algorithms to solve this problem, and we study several cases. If there is no bandwidth constraint, we propose a polynomial algorithm that compute the optimal path. We also give lower and upper bounds on the optimal path length. On the other hand, we show that the problem is NP-hard if there is a bandwidth constraint (or other Quality of Service parameters), even if there is only two protocols and in a symmetric graph. We study the complexity and the scalability of our algorithms and evaluate their performances on real and random topologies. The results show that they are faster than the previous ones proposed in the literature. These algorithms can also have important applications in automatic tunneling.

I. INTRODUCTION

Carrier-grade networks generally encompass several layers that involve different technologies and protocols. To support some services, such as a Virtual Private Network (VPN), a path across network equipments must be identified and the equipments be configured accordingly. Under stringent requirements of Quality of Service (QoS – e.g., end-to-end delay, geographic zone avoidance, etc.), computing such a path within a single layer is not always possible. Hence, one of the key challenges is to determine the end-to-end path that uses the appropriate *adaptation functions* over the protocols: The mapping from a protocol to another being realized through *encapsulation* (e.g., Ethernet over IP/MPLS [1]), *decapsulation* (the reverse operation) or *conversion* (e.g., IPv4 to IPv6 [2]) functions. Consequently, the path computation process should take into account the adaptation function capabilities of the network equipments in order to ensure *path feasibility*: If a protocol is encapsulated in another one, then it must be decapsulated (or unwrapped) further in the path. If several encapsulations are nested, then the corresponding decapsulations must occur in

the right order. Here, the multi-layer context should be taken in a broad sense: Presence of several protocols and technologies that can be nested, encapsulated, converted, etc.

Dealing with protocol heterogeneity is increasingly important nowadays. In addition to the IPv4/IPv6 migration, this heterogeneity appears in tunneling, some architectures (e.g., The Pseudo-Wire architecture [3] allows the emulation – and thus the encapsulation – of lower layer protocols over Packet-Switched Networks), hybrid networks (e.g., National Research and Education Networks – NRENs – which may have optical and IP interconnection points), and last but not least, most carrier-grade networks, which have separate control planes for IP and Transport layers. In all these contexts, a unified control plane would be very useful for optimizing the network resources and for reducing operational and management costs.

OpenFlow, and specifically protocol-agnostic versions such as Protocol-Oblivious Forwarding (POF) [4] or P4 [5], are an opportunity to design such a control plane. Some previous works [6], [7] present an OpenFlow-based architecture to achieve this challenge, but they only focus on the convergence of packet and circuit networks. Other works tackle the traffic engineering problem in SDNs but they circumscribe it to a single layer [8] or to the IPv4/IPv6 migration [9]. However, an important problem to solve remains the path computation process in a multi-layer context. Taking into account the adaptation functions is not trivial and classical algorithms such as (Edsger) Dijkstra’s one cannot achieve the task as they do not handle these functions. Here, we design several algorithms to compute minimum cost paths dealing with protocol changes and adaptation functions.

Our contributions:

- 1) We widely generalize the model and the polynomial algorithms described in Lamali *et al.* [10] to perform path computation in multi-layer networks (without bandwidth constraint). Our model takes into account all possible types of protocol changes (encapsulation, conversion, etc.) and any additive metric. We drastically improve the algorithm complexity and show its efficiency through simulations on real and random topologies.
- 2) For simulation purposes, we empirically study the distribution of adaptation functions over the network nodes and its impact on feasible path existence. We exhibit a phase transition phenomenon, i.e., a gap where the probability of existence of a feasible path hugely increases.

- 3) We provide lower and upper bounds on the length of the shortest path: We show that the shortest path length can be superpolynomial in the network size, but cannot be more than doubly exponential.
- 4) We prove that path computation in multi-layer networks under bandwidth constraint is NP-hard even with two protocols and on symmetric graphs, thus improving a result of Kuipers and Dijkstra [11]. We also obtain results on the complexity of some subproblems: It is polynomial on *Directed Acyclic Graphs* (DAG) and the general problem is not approximable. We propose a new heuristic to resolve the problem and we show its efficiency through simulations.
- 5) We propose the first algorithm to perform path computation in multi-layer networks under several QoS constraints by adapting the Self-Adaptive Multiple Constraints Routing Algorithm (SAMCRA – Van Mieghem and Kuipers [12]) to the multi-layer context. We study its scalability through simulations.

The most important possible applications of our algorithms are the unification of control planes and automatic (nested) tunneling. This paper extends the work published in [13]. It provides the detailed algorithms, together with their proofs of correctness and complexity analysis. It also provides lower and upper bounds on the shortest path length. This has an important impact on the complexity of the problem.

The paper is organized as follows: Section II describes the problem of path computation in multi-layer networks, and recalls the related work; Section III formalizes the problem and describes our model of multi-layer network; Section IV proposes algorithms to perform path computation without bandwidth constraint, and shows their efficiency through simulations, it also studies the phase transition phenomenon in multi-layer networks; Section V provides lower and upper bounds on the shortest path length; Section VI studies the complexity of path computation under bandwidth constraint and proposes heuristic solutions to tackle the problem; Section VII proposes the first algorithm computing paths under additive QoS constraints and studies its scalability; finally, Section VIII concludes the paper.

II. PATH COMPUTATION IN MULTI-LAYER NETWORKS

A. Connectivity in multi-layer networks

We aim to present the different concepts of path computation in multi-layer networks through an example. While this example relates to multi-domain multi-layer networks, the underlying problem of path computation is the same as in a single domain network¹. Figure 1 (inspired by the Inter-Provider Reference Model [15]) depicts a network involving multiple domains and adaptation function capabilities of network equipments: A company owning a Local Area Network (LAN) wishes the Virtual Machines (VMs) of a data-center to be within the same routing domain (for instance through

a Layer 2 VPN or a Generic Routing Encapsulation tunnel). Hence, the switches of the LAN and the VMs of the data-center must communicate through Ethernet datagrams and a path has to be determined across Domains 1 and 2.

On Figure 1, Domains 1 and 2 use IPv6/MPLS-TE technology and are linked by equipments providing Ethernet encapsulation and decapsulation. The Provider Edge (PE) of Domain 1 is linked to the Customer Edge (CE) of the data-center. The adaptation capabilities of each node are shown above it. An instance of a feasible path would cross the PE of Domain 1 converting IPv4 packets into IPv6 ones. Then it would apply the encapsulation and decapsulation of the border routers of Domains 1 and 2 respectively, and the PE of Domain 2 would apply a conversion of IPv6 packets into IPv4 ones. The protocol stacks of the packets at each stage are illustrated at the bottom of Figure 1. As an example of unfeasible path, a direct Ethernet connection between the CE of the data-center and the border router of Domain 1 appears. This configuration leads to a decapsulation of an IPv6 packet from an Ethernet datagram (by the border router of Domain 2) whereas at this stage the datagram encapsulates IPv4 packets.

This example depicts the constraints to comply with when computing a multi-layer (and multi-domain, in this case) path: Being physically linked is not sufficient to establish connectivity. Protocol continuity (by analogy with wavelength continuity in optical networks) must hold and the adaptation functions should occur in the right order. Moreover, feasible paths can involve loops and their subpaths are not necessary feasible [11], [16]. Nowadays, such paths are manually determined and configurations are operated and applied by scripts.

B. Related work

The initial works dealing with protocol and technology heterogeneity circumscribed the problem at the optical layer. For instance, Chlamtac *et al.* [17] described a model and algorithms to compute a path under wavelength continuity constraints. Zhu *et al.* [18] addressed the same problem in WDM mesh networks tackling traffic grooming issues. In [19], Gong and Jabbari provided an algorithm to compute an optimal path under constraints on several layers: wavelength continuity, label continuity, etc.

However, the models of these past works are not adapted to the problem of nested encapsulation and decapsulation capabilities, for which a kind of stack mechanism is needed. In [20], Dijkstra *et al.* addressed this issue in the context of the ITU-T G.805 recommendations on adaptation functions. They stressed the lack of solutions on path computation. Kuipers and Dijkstra [11] demonstrated that the problem of path computation with encapsulation and decapsulation capabilities is NP-hard under bandwidth constraint. They proposed a Breadth-First Search (BFS) algorithm that explores all possible paths until finding a feasible one. In [10], Lamali *et al.* proved that the problem is polynomial if the bandwidth constraint is relaxed. Their approach was to model the network as a Push-Down Automaton and to use automata and language theory tools to compute a shortest feasible path, but only considering the number of hops or adaptation functions. More recently,

¹The algorithms presented in this paper can be applied in a single-domain or a multi-domain context. For the latter, however, a mechanism for sharing the network information (such as the topology) is needed. This can be done through a PCE for example [14].

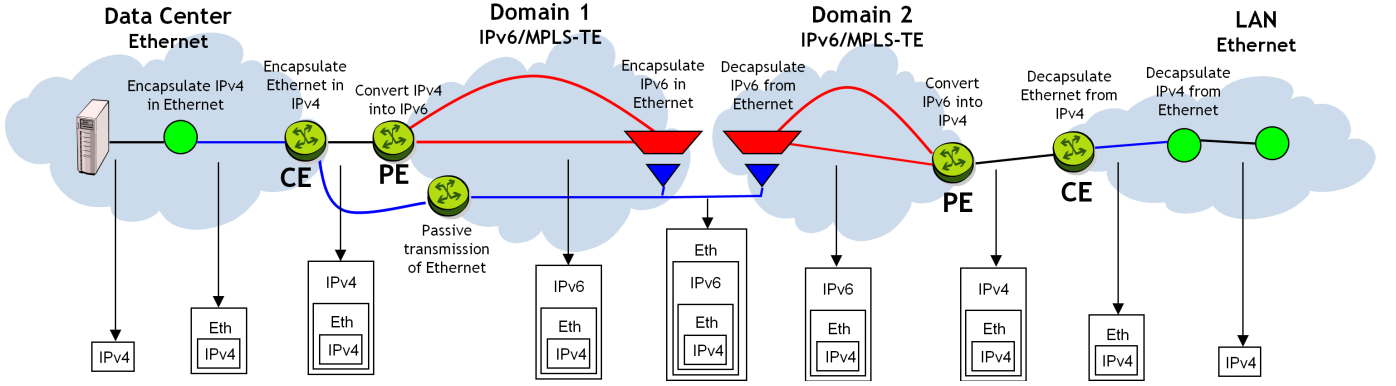


Fig. 1. Carrier-grade network comprising several domains and different layers.

Iqbal *et al.* [21] underlined the need of path computation algorithms in NRENs. They proposed a new matrix-based model for multi-layer networks and algorithms based on k -shortest paths and *LOOK-AHEAD* methods. However, the model deals with technologies² instead of protocols. Thus, the nested protocols are not transparent to the nodes. Moreover, the proposed exact algorithm is exponential and can compute only loopless feasible paths.

C. Proposed approach

Our goal is to study the path computation problem in a multi-layer context and to propose efficient algorithms to resolve it. To this end, we focus on three cases: Path computation without bandwidth constraint (by adapting the language theoretic approach of Lamali *et al.* [10]), under bandwidth constraint (by using graph transformation in order to overcome the problem complexity) and under several QoS constraints. The simulations showing the efficiency of our algorithms follow a methodology based on the probabilistic distribution of the adaptation functions over the nodes.

III. MODEL AND PROBLEM FORMALIZATION

This section describes a mathematical model of multi-layer networks and formalizes the notion of path feasibility.

A. Multi-layer network model

Notation convention. In order to avoid confusion, lowercase letters denote protocols (e.g., a, b, c, x, y) or functions (e.g., f, h, ℓ). Capital letters denote nodes and links (e.g., U, V, E). Finally, calligraphic letters denote sets (e.g., $\mathcal{G}, \mathcal{V}, \mathcal{E}$).

A multi-layer network is a 4-tuple $\mathcal{N} = (\mathcal{G}, \mathcal{A}, \mathcal{F}, h)$ where:

- $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ is a directed graph modeling the network topology. The set of nodes \mathcal{V} models the routers of the network. The set of edges \mathcal{E} models the physical links between the routers. The number of nodes is denoted by $|\mathcal{V}| = n$ and the number of links is denoted by $|\mathcal{E}| = m$.
- $\mathcal{A} = \{a, b, c, \dots\}$ is the set of protocols available in the network, but not necessarily at each router. The number of protocols is denoted by $|\mathcal{A}| = \lambda$.

²A *technology* is an exhaustive description of the protocol stack at some node, e.g., IP over Ethernet over ATM.

- For each node $U \in \mathcal{V}$, $\mathcal{F}(U)$ is the set of adaptation functions available on node U . These functions are:
 - *Conversion*: A protocol a is converted into a protocol b without any change of the possible underlying protocols. It is denoted by $(a \rightarrow b)$. E.g., Wavelength conversion on the optical layer, IPv4 to IPv6, etc.
 - *Passive function*: A protocol a is left as it is. It is a classical retransmission without any protocol change and can be considered as a special case of protocol conversion where $a = b$. It is denoted by $(a \rightarrow a)$.
 - *Encapsulation*: A protocol a is encapsulated in a protocol b . It is denoted by $(a \rightarrow ab)$.
 - *Decapsulation*: A protocol a is decapsulated from a protocol b . It is denoted by $(a \rightarrow \overline{ab})$.
- $h : \mathcal{V} \times \mathcal{F} \times \mathcal{V} \rightarrow \mathbb{R}_+$ is the *weight* function. The value $h(U, f, V)$ (where $U, V \in \mathcal{V}$ and $f \in \mathcal{F}(U)$) is the cost of using the link (U, V) with the adaptation function f on U . Hence, function h allows representing any additive metric either associated only to the links or to both links and adaptation functions.

B. Path feasibility

Let (S, D) be a pair of nodes in \mathcal{G} corresponding to the source and the destination of the path to be computed. We consider a path from S to D as a sequence of nodes and adaptation functions $S f_0 U_1 f_1 U_2 f_2 \dots U_k f_k D$ where each U_i , $i = 1, \dots, k$, is a node and each f_i is an adaptation function (f_0 being dummy). A path is *feasible* if:

- 1) The sequence $S U_1 U_2 \dots U_k D$ is a path in $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ and each $f_i \in \mathcal{F}(U_i)$;
- 2) Each encapsulated protocol is decapsulated before reaching D according to its encapsulation order and *protocol continuity* holds (i.e., if the sequence contains a function f_i s.t. $f_i = (a \rightarrow b)$, $a, b \in \mathcal{A}$, then $f_{i+1} = (b \rightarrow a')$ or $f_{i+1} = (b \rightarrow ba')$ or $f_{i+1} = (\overline{a'} \rightarrow a'b)$, $a' \in \mathcal{A}$).

Actually, the protocol sequences of feasible paths can be characterized as a well-parenthesized language [10].

IV. PATH COMPUTATION WITHOUT BANDWIDTH CONSTRAINT

This section proposes a polynomial algorithm to resolve the path computation problem without bandwidth constraint and

evaluates it through simulations. Lamali *et al.* [10] proposed a language theoretic approach to compute a shortest feasible path (involving encapsulations and decapsulations of protocols) in a multi-layer network. The metric considered was the number of hops or of encapsulations in the path. The approach comprises the following steps:

- 1) Consider the set of protocols as an alphabet and the multi-layer network as a Push-Down Automaton (PDA);
- 2) If the considered metric is the number of encapsulations, transform the PDA in order to bypass passive transitions;
- 3) Convert the PDA to a Context-Free Grammar (CFG);
- 4) Compute the shortest word generated by the CFG. It is the protocol sequence of a shortest path;
- 5) Compute a shortest path from this sequence.

We made several improvements to these algorithms:

- The PDA building is modified in order to support protocol conversion by adding a new transition type;
- The PDA transitions are weighted in order to reflect the weight function. Thus, our algorithm computes the minimum cost path according to any additive metric (instead of just the number of hops or encapsulations);
- The PDA transformation is no longer useful thanks to the weight function: Simply put $h(U, f, V) = 1$ (where $U, V \in \mathcal{V}$ and $f \in \mathcal{F}(U)$) for all triples where f is an encapsulation, and $h(U, f, V) = 0$ for all other triples.
- The conversion of the PDA into a CFG is adapted in order to assign weights to the production rules;
- Since the production rules are weighted, the goal is no longer to compute the shortest protocol sequence but the one having the minimum weight derivation tree. This is done thanks to Knuth's algorithm described in [22].
- The computation of the path matching the protocol sequence is modified to take into account the weights.

Additionally to these improvements, the algorithm complexity is drastically decreased from $O(\lambda^8 n^7)$ in [10] to $O(\lambda^5 n^2 m)$.

A. Detailed algorithms

1) Theoretical language aspects of multi-layer paths

Considering a path $\mathcal{P} = S f_0 U_1 f_1 U_2 f_2 \dots U_k f_k D$, let $H_{\mathcal{P}} = f_1 \dots f_k$ denotes the sequence of adaptation functions along \mathcal{P} . Let define as an alphabet the set $\overline{\mathcal{A}} = \{\overline{a} \mid a \in \mathcal{A}\}$ and the set $\underline{\mathcal{A}} = \{\underline{a} \mid a \in \mathcal{A}\}$.

$\mathcal{T}_{\mathcal{P}} = x_1 \dots x_{k+1}$ is the sequence of protocols used along path \mathcal{P} . It is called the *trace* of \mathcal{P} . For each x_i :

- $x_i = a$ and $x_{i+1} = b, \overline{b}$ or \underline{b} means that U_i converts protocol a into b ($a, b, \overline{b}, \underline{b} \in \mathcal{A} \cup \overline{\mathcal{A}} \cup \underline{\mathcal{A}}$)
- $x_i = \overline{a}$ and $x_{i+1} = b, \overline{b}$ or \underline{b} means that U_i encapsulates protocol a in b
- $x_i = \underline{a}$ and $x_{i+1} = b, \overline{b}$ or \underline{b} means that U_i decapsulates protocol b from a .

Here, some additional definitions are needed. The set of protocol conversions available on node U is denoted by $\mathcal{CO}(U)$. The set of encapsulations available on node U is denoted by $\mathcal{EN}(U)$ and the set of decapsulations available on node U is denoted by $\mathcal{CD}(U)$.

$In(U)$ (resp. $Out(U)$) is the set of protocols that node U can receive (resp. send). More formally:

- If $(a \rightarrow b) \in \mathcal{CO}(U)$ then $a \in In(U)$ and $b \in Out(U)$
- If $(a \rightarrow ab) \in \mathcal{EN}(U)$ then $a \in In(U)$ and $b \in Out(U)$
- If $(a \rightarrow ab) \in \mathcal{DE}(U)$ then $b \in In(U)$ and $a \in Out(U)$

Several paths can have the same trace. The set of traces of the feasible paths in a network \mathcal{N} is a context-free language but it is not regular as the encapsulations and decapsulations should be balanced. In fact, it is a context-free language, and thus requires a stack to be recognized and computed. PDAs are the classical tools to recognize context-free languages. Using weighted PDAs allows associating a weight to each link and adaptation function in order to model any additive metric.

2) Definition of WPDA

A weighted PDA (WPDA) is a 8-tuple $\mathcal{PD}\mathcal{A} = (\mathcal{S}, \Sigma, \Gamma, \delta, Q_0, Z_0, \mathcal{S}_F, \omega)$ where \mathcal{S} is the set of states, Σ is the input alphabet, Γ is the stack symbol set (i.e., stack alphabet) not necessarily different from Σ , δ is the set of transitions, Q_0 is the initial state, Z_0 is the initial stack symbol, \mathcal{S}_F is the set of final (accepting) states and ω is a weight function over the set of transitions (i.e., $\omega : \delta \rightarrow \mathbb{R}_+$). A transition $t \in \delta$ is denoted by $t = (Q_i, \langle x, \alpha, \beta \rangle, Q_j)$, where Q_i (resp. Q_j) is the state of $\mathcal{PD}\mathcal{A}$ before (resp. after) the transition, $x \in \Sigma \cup \{\epsilon\}$ is an input symbol, $\alpha \in \Gamma$ is the symbol that is popped from the top of the stack, and $\beta \in \Gamma^*$ is the symbol sequence that is pushed on the top of the stack.

Remark. WPDAs are more often formalized as 6-tuples $(\mathcal{S}, \Gamma, \mathcal{M}, Q_0, Z_0, \mathcal{S}_F)$ where \mathcal{M} , called the *Push-Down transition matrix*, is a matrix over a semiring of formal power series. The input alphabet Σ , the transitions set δ and the weight function ω are expressed by a single entity $\mathcal{M} \in ((\mathcal{R}\langle\langle\Sigma^*\rangle\rangle)^{\mathcal{S} \times \mathcal{S}})^{\Gamma^* \times \Gamma^*}$, where $\mathcal{R}\langle\langle\Sigma^*\rangle\rangle$ denotes the collection of all power series from Σ^* into a semiring \mathcal{R} . For simplicity, we opt for defining a WPDA as a classical PDA with a weight function over the transition set. For the theoretical foundations of WPDAs, the interested reader can refer to [23].

3) From the graph to the WPDA

Algorithm 1 converts a multi-layer network \mathcal{N} with a specified pair of nodes (S, D) into a WPDA $\mathcal{PD}\mathcal{A} = (\mathcal{S}, \Sigma, \Gamma, \delta, Q_0, Z_0, \mathcal{S}_F = \{Q_F\}, \omega)$.

Computing a feasible path requires to know the current protocol and the last encapsulated one (in order to know if a decapsulation can be performed). Thus Algorithm 1 creates a state U_x for each node U and each protocol $x \in In(U)$. Being in a state U_x indicates that the current protocol is x . The last encapsulated protocol is the one on the top of the stack. The conversion functions $(x \rightarrow y)$ between node U and node V are turned into transitions $(U_x, \langle x, \alpha, \alpha \rangle, V_y)$ in the WPDA. The encapsulation functions $(x \rightarrow xy)$ are converted into pushes of x on the stack $(U_x, \langle \overline{x}, \alpha, x\alpha \rangle, V_y)$ and the decapsulation functions into pops of x from the stack $(U_y, \langle \underline{y}, x, \emptyset \rangle, V_x)$.

Complexity of Algorithm 1. The complexity of Algorithm 1 is in $O(\lambda^3 m)$. The number of states created is at worst $2 + \lambda(n - 1)$. The number of transitions created is in $O(\lambda^3 m)$, which is also the complexity of the whole algorithm.

Proposition 1. A path \mathcal{P} in a network \mathcal{N} is feasible if and only if its trace $\mathcal{T}_{\mathcal{P}}$ is accepted by $\mathcal{PD}\mathcal{A}$.

Proof. Let $\mathcal{P} = S f_0 U_1 f_1 U_2 f_2 \dots U_k f_k D$ be a feasible path. By construction, for each 3-tuple (U_i, f_i, U_{i+1}) there is a

Algorithm 1 Convert a network into a WPDA

Input: A network $\mathcal{N} = (\mathcal{G} = (\mathcal{V}, \mathcal{E}), \mathcal{A}, \mathcal{F}, h)$, a source S and a destination D

Output: A WPDA $\mathcal{PDA} = (S, \Sigma, \Gamma, \delta, Q_0, Z_0, \{Q_F\}, \omega)$
 $\Sigma \leftarrow \mathcal{A} \cup \overline{\mathcal{A}} \cup \underline{\mathcal{A}}; \Gamma \leftarrow \mathcal{A}$
 Create a single state Q_0 corresponding to node S
 Create a dummy final state Q_F
for each node $U \neq S$ in \mathcal{V} and each $x \in In(U)$ **do**
 Create a state U_x
for each state U_x s.t. $(S, U) \in \mathcal{E}$ and each $x \in Out(S)$ **do**
 Create the transition $t = (Q_0, \langle \epsilon, Z_0, Z_0 \rangle, U_x)$
 $\omega(t) \leftarrow 0$
for each $x \in In(D)$ **do**
 Create the transition $t = (D_x, \langle x, Z_0, \emptyset \rangle, Q_F)$
 $\omega(t) \leftarrow 0$
for each link $(U, V) \in \mathcal{E}$ s.t. $U \neq S$ **do**
 for each $(x \rightarrow y) \in \mathcal{CO}(U)$ **do**
 if $y \in In(V)$ **then**
 for all $\alpha \in \Gamma$ **do**
 Create the transition $t = (U_x, \langle x, \alpha, \alpha \rangle, V_y)$
 $\omega(t) \leftarrow h(U, (x \rightarrow y), V)$
 for each $(x \rightarrow xy) \in \mathcal{EN}(U)$ **do**
 if $y \in In(V)$ **then**
 for all $\alpha \in \Gamma$ **do**
 Create the transition $t = (U_x, \langle \bar{x}, \alpha, x\alpha \rangle, V_y)$
 $\omega(t) \leftarrow h(U, (x \rightarrow xy), V)$
 for each $(x \rightarrow \underline{xy}) \in \mathcal{DE}(U)$ **do**
 if $x \in In(V)$ **then**
 Create the transition $t = (U_y, \langle \underline{y}, x, \emptyset \rangle, V_x)$
 $\omega(t) \leftarrow h(U, (x \rightarrow \underline{xy}), V)$

transition:

- $t = ((U_i)_x, \langle x, \alpha, \alpha \rangle, (U_{i+1})_y)$ if $f_i = (x \rightarrow y)$
- $t = ((U_i)_x, \langle \bar{x}, \alpha, x\alpha \rangle, (U_{i+1})_y)$ if $f_i = (x \rightarrow xy)$
- $t = ((U_i)_x, \langle \underline{x}, y, \emptyset \rangle, (U_{i+1})_y)$ if $f_i = (\overline{y \rightarrow yx})$

This transition recognizes the i^{th} letter of the trace $\mathcal{T}_{\mathcal{P}}$. It is easy to show by induction that $\mathcal{T}_{\mathcal{P}}$ is accepted by \mathcal{PDA} .

Conversely, if a trace $\mathcal{T}_{\mathcal{P}}$ is accepted by a transition sequence $t_1 \dots t_k$ where each $t_i = ((U_i)_x, \langle x, \alpha, \beta \rangle, (U_{i+1})_y)$. Then there is an adaptation function:

- $f_i = (x \rightarrow y) \in \mathcal{CO}(U_i)$ if the transition $t_i = ((U_i)_x, \langle x, \alpha, \alpha \rangle, (U_{i+1})_y)$
- $f_i = (x \rightarrow xy) \in \mathcal{EN}(U_i)$ if the transition $t_i = ((U_i)_x, \langle \bar{x}, \alpha, x\alpha \rangle, (U_{i+1})_y)$
- $f_i = (\overline{y \rightarrow yx}) \in \mathcal{DE}(U_i)$ if the transition $t_i = ((U_i)_x, \langle \underline{x}, y, \emptyset \rangle, (U_{i+1})_y)$

Thus the path $Sf_0U_1f_1U_2f_2 \dots U_kf_kD$ is feasible in \mathcal{N} . \square

The weight of a path $\mathcal{P} = Sf_0U_1f_1U_2f_2 \dots U_kf_kD$ is the sum of the weights of its links and its adaptation functions. It is denoted by $h(\mathcal{P}) \stackrel{def}{=} \sum_{i=1}^k h(U_i, f_i, U_{i+1})$ with $U_{k+1} = D$. The weight of a transition sequence is the sum of the weights of each transition (i.e., $\omega(\{t_1, t_2, \dots, t_k\}) = \sum_{i=1}^k \omega(t_i)$). The weight of a word w , denoted by $\omega(w)$, is the weight of the transitions that accept w in \mathcal{PDA} . But as \mathcal{PDA}

may be nondeterministic, a same word may be accepted by different transition sequences. We consider only the transition sequence of minimum weight that accepts w . More formally, $\omega(w) = \min_{t_1, \dots, t_k \in \delta} \omega(\{t_1 \dots t_k\})$ s.t. $\{t_1 \dots t_k\}$ accepts w .

Lemma 1. *If \mathcal{PDA} accepts the trace $\mathcal{T}_{\mathcal{P}}$ of a path \mathcal{P} , then $\omega(\mathcal{T}_{\mathcal{P}}) = h(\mathcal{P}^*)$, where \mathcal{P}^* is the path of minimum cost having $\mathcal{T}_{\mathcal{P}}$ as trace.*

Proof. By definition, $\omega(\mathcal{T}_{\mathcal{P}}) = \omega(\{t_1, \dots, t_k\})$, where $\{t_1, \dots, t_k\}$ is the transition sequence with minimal weight which accepts $\mathcal{T}_{\mathcal{P}}$. From $\{t_1, \dots, t_k\}$, it is possible to build the path \mathcal{P}^* (reversing the conversion in Algorithm 1) such that $\mathcal{T}_{\mathcal{P}^*} = \mathcal{T}_{\mathcal{P}}$, and $h(\mathcal{P}^*) = \omega(\mathcal{T}_{\mathcal{P}})$. Suppose that $\exists \mathcal{P}'$ s.t. $\mathcal{T}_{\mathcal{P}'} = \mathcal{T}_{\mathcal{P}}$ and $h(\mathcal{P}') < \omega(\mathcal{T}_{\mathcal{P}})$, then it is possible to build from \mathcal{P}' a sequence of transitions that corresponds to the links and adaptation functions involved in \mathcal{P}' (as in Algorithm 1). Let this sequence be $t'_1 \dots, t'_k$. The weight of each transition t'_i corresponds to the weight of an adaptation function associated to a link in \mathcal{P}' . The weight of $t'_1 \dots, t'_k$ is then less than $\omega(\mathcal{T}_{\mathcal{P}})$, and by Proposition 1, this sequence accepts $\mathcal{T}_{\mathcal{P}}$. This is inconsistent with the definition of $\omega(\mathcal{T}_{\mathcal{P}})$. \square

4) From the WPDA to a WCFG

In order to compute the minimum weight trace and its corresponding path, \mathcal{PDA} is converted into a weighted Context-Free Grammar (WCFG). A WCFG is a CFG with a weight function over the set of production rules. The conversion of a PDA into a CFG is well-known. The conversion of a WPDA into a WCFG is done in the same way, in addition the weight of each transition is assigned to the corresponding production rules (called rules in Algorithm 2) in the WCFG. Algorithm 2 is an adaptation of the method described in [24]. It converts \mathcal{PDA} into a WCFG $\mathcal{WCFG} = (\mathcal{Q}, \Sigma, [Q_0], \mathcal{R}, \pi)$ where:

- \mathcal{Q} is the set of nonterminals,
- Σ is the alphabet or set of terminals (the same as the WPDA input alphabet),
- $[Q_0]$ is the initial symbol (initial nonterminal, or axiom),
- \mathcal{R} is the set of production rules,
- $\pi : \mathcal{R} \rightarrow \mathbb{R}_+$ is the weight function over the set of production rules.

Complexity of Algorithm 2. The number of nonterminals is bounded by $O(|\Gamma| \times |\mathcal{S}|^2)$. The number of production rules is bounded by $O(|\delta| \times |\mathcal{S}|^2)$, which is the worst case complexity of Algorithm 2. This corresponds to $O(\lambda^5 n^2 m)$.

The minimum weight derivation tree: Generating the minimum weight trace (and then the minimum weight path) requires to build its derivation tree. Let $[X]$ be a nonterminal, we define $\ell([X])$ as the sum of the weights of the productions needed for, starting from $[X]$, deriving a word in Σ^* . Thus $\ell([Q_0])$ is the weight of the minimum weight trace.

The function is $\ell : \{\mathcal{Q} \cup \Sigma \cup \{\epsilon\}\}^* \rightarrow \mathbb{N} \cup \{\infty\}$ s.t.:

- if $w = \epsilon$ or $w \in \Sigma$ then $\ell(w) = 0$,
- if $w = \alpha_1 \dots \alpha_k$ (with $\alpha_i \in \{\mathcal{Q} \cup \Sigma \cup \{\epsilon\}\}$) then $\ell(w) = \sum_{i=1}^k \ell(\alpha_i)$.
- Let $r_1 = [X] \rightarrow \gamma_1, r_2 = [X] \rightarrow \gamma_2, \dots, r_k = [X] \rightarrow \gamma_k$ be the set of production rules having $[X]$ as left part. Then $\ell([X]) = \min\{\pi(r_1) + \ell(\gamma_1), \dots, \pi(r_k) + \ell(\gamma_k)\}$

Algorithm 2 Convert a WPDA into a WCFG

Input: $\mathcal{PDA} = (\mathcal{S}, \Sigma, \Gamma, \delta, Q_0, Z_0, \{Q_F\}, \omega)$
Output: $\mathcal{CFG} = (\mathcal{Q}, \Sigma, [Q_0], \mathcal{R}, \pi)$
 Create the axiom $[Q_0]$
for each state $U_x \in \mathcal{S}$ **do**
 Create nonterminal $[Q_0 Z_0 U_x]$ and rule $[Q_0] \rightarrow [Q_0 Z_0 U_x]$
for each transition $(U_x, \langle x, \alpha, \beta \rangle, V_y)$ **do**
 if $\beta = \emptyset$ (pop) **then**
 Create a nonterminal $[U_x \alpha V_y]$ and rule $[U_x \alpha V_y] \rightarrow x$
 $\pi(r) \leftarrow \omega(U_x, \langle x, \alpha, \emptyset \rangle, V_y)$
 if $\beta = \alpha$ (conversion transition) **then**
 for each $Q_i \in \mathcal{S}$ **do**
 Create nonterminals $[U_x \alpha Q_i]$ and $[V_y \alpha Q_i]$
 Create the rule $r = [U_x \alpha Q_i] \rightarrow x[V_y \alpha Q_i]$
 $\pi(r) \leftarrow \omega(U_x, \langle x, \alpha, \beta \rangle, V_y)$
 if $\beta = x\alpha$, $x \in \Gamma$ (push) **then**
 for each $(Q_i, Q_j) \in \mathcal{S}^2$ **do**
 Create nonterminals $[U_x \alpha Q_j]$, $[V_y \alpha Q_i]$ and $[Q_i \alpha Q_j]$
 Create the rule $r = [U_x \alpha Q_j] \rightarrow x[V_y \alpha Q_i][Q_i \alpha Q_j]$
 $\pi(r) \leftarrow \omega(U_x, \langle x, \alpha, x\alpha \rangle, V_y)$

Knuth's algorithm [22] can be adapted to compute the minimum weight derivation tree of a grammar. This corresponds to the weight of \mathcal{T}_P , where P is the minimum cost path to compute. The adapted algorithm maintains a list of production rules and updates the $\ell[X]$ according to the formula above. The sketch of the algorithm is as follows:

- Initialize $\ell([X])$ to ∞ for each nonterminal $[X]$
- For each production rule $[X] \rightarrow \alpha_1 \dots \alpha_n$:
 $\ell([X]) \leftarrow \min\{\ell([X]), \pi(r) + \sum_{i=1}^n \ell(\alpha_i)\}$

When the value of $\ell([X])$ is updated, the production rule that originates this update is stored as in the classical Dijkstra's algorithm. This allows to generate the desired word.

The algorithm terminates when all the $\ell[X]$ have the right value and no additional update is possible. Implementing this algorithm with Fibonacci heaps leads to a $O(|\mathcal{Q}| \log |\mathcal{Q}| + |\mathcal{R}|)$ complexity [25], which corresponds to $O(\lambda^5 n^2 m)$.

5) Deriving the minimum cost path from its trace

Algorithm 3 generalizes the algorithm proposed in [10]. It takes as input the minimum weight trace \mathcal{T}_P accepted by \mathcal{PDA} and computes the path \mathcal{P} that matches it³.

Algorithm 3 starts on $nodes[1] = S$ then checks at each step all the links in \mathcal{E} that match the current letter (protocol) in \mathcal{T}_P . If $\mathcal{T}_P = x_1 x_2 \dots x_k$ ($x_i \in \mathcal{A} \cup \bar{\mathcal{A}} \cup \underline{\mathcal{A}}$), then at each step i , the algorithm starts from each node U in $nodes[i]$ and adds to $links[i]$ all the links (U, V) which match x_i . Each V is added in $nodes[i+1]$. The value $weights[(U, V), i]$ is the cost of using link (U, V) at step i . It corresponds to the weight $h(U, f_i, V)$ where f_i is the adaptation function used at step i . When the trace \mathcal{T}_P is covered, a classical minimum cost path algorithm from S to D in the graph $(nodes, links, weights)$ computes the minimum weight path.

Complexity of Algorithm 3. The complexity of Algorithm 3

Algorithm 3 Computing the minimum cost path

Input: The network \mathcal{N} and \mathcal{T}_P
Output: The minimum cost path \mathcal{P}
 $nodes[1] \leftarrow S$; $i \leftarrow 2$
while The trace is not completely covered **do**
 for each $U \in nodes[i]$, $V \in \mathcal{V}$ s.t. $(U, V) \in \mathcal{E}$ **do**
 if $x_i \in \mathcal{A}$, $x_i \in Out(U)$, $x_i \in In(V)$ and $(x_{i-1} \rightarrow x_i) \in \mathcal{CO}(U)$ **then**
 Add (U, V) in $links[i]$ and V in $nodes[i+1]$
 $weights[(U, V), i] \leftarrow h(U, (x_{i-1} \rightarrow x_i), V)$
 if $x_i \in \bar{\mathcal{A}}$, $x_i \in Out(U)$, $x_i \in In(V)$ and $(x_{i-1} \rightarrow x_{i-1}x_i) \in \mathcal{EN}(U)$ **then**
 Add (U, V) in $links[i]$ and V in $nodes[i+1]$
 $weights[(U, V), i] \leftarrow h(U, (x_{i-1} \rightarrow x_{i-1}x_i), V)$
 if $x_i \in \underline{\mathcal{A}}$, $x_i \in Out(U)$, $x_i \in In(V)$ and $(x_i \rightarrow x_i x_{i-1}) \in \mathcal{DE}(U)$ **then**
 Add (U, V) in $links[i]$ and V in $nodes[i+1]$
 $weights[(U, V), i] \leftarrow h(U, (x_i \rightarrow x_i x_{i-1}), V)$
 $i++$
 Compute the minimum cost path from S to D in $(nodes, links, weights)$

is bounded by $O(|\mathcal{T}_P|nm)$ in the worst case.

B. Simulations

We implemented our algorithm (called PDA) and compared it to a classical BFS approach.

1) Networks used for the simulations and methodology

Large multi-layer topologies are generally not available. Some public ones as the Internet2 network [26] are not large enough to show the scaling of our algorithm. Thus, we performed simulations on two real topologies described in [27] and random scale-free topologies:

- Topology $T1$ is a simplified version of Time Warner network. It has 41 nodes and 296 directed links.
- Topology $T2$ corresponds to the network of Exodus as in 2002. It has 79 nodes and 294 directed links.
- Scale-free topologies are generated according to the Barabasi-Albert model [28]. A complete graph of 10 nodes is created. Then nodes are iteratively added according to a preferential attachment mechanism: A new node is randomly linked to 5 existing nodes. The probability for a node to be selected is proportional to its degree.

Since these topologies are not layered, the adaptation functions are randomly allocated to the nodes. For a set of protocols \mathcal{A} , there are $3\lambda^2$ possible adaptation functions (for each ordered pair of protocols: a conversion, an encapsulation and a decapsulation). For each node U , each of these adaptation functions is available on U with probability p . The source and the destination nodes are the diameter extremities, which corresponds to 5 (resp. 10) hops for Topology $T1$ (resp. $T2$).

2) Phase transition in path feasibility

Depending on the network topology and the adaptation function distribution, there is not always a feasible path. It is interesting to know the probability of a feasible path

³It is possible that several paths match the trace. In this case, the path can be chosen randomly or according to a load-balancing policy.

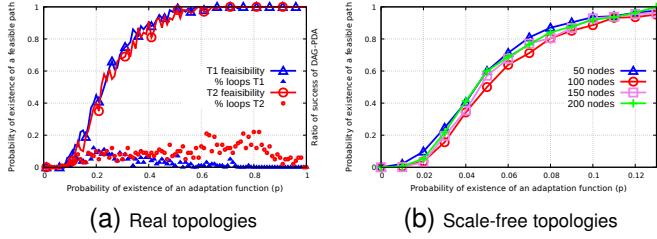


Fig. 2. Probability of existence of a feasible path in real and scale-free topologies according to the probability of existence of an adaptation function.

existence according to probability p in order to set appropriate parameters for the simulations. In case of path existence, knowing the probability that the minimum cost one involves loops allows comparing the different algorithms (some of them allow loops and others do not). To compute this probability, we performed 200 runs for each value of p and counted the number of times there was a feasible path.

Figure 2a shows the evolution of feasible path existence probability according to p and the proportion of minimum cost paths that involve loops. Not surprisingly, the probability of feasible path existence grows according to p . On both topologies, the probability of path existence reaches 50% when $p = 0.22$ and follows a phase transition phenomenon. For example, in the interval $p \in [0.10, 0.38]$, the probability of path existence in Topology $T1$ grows from 5% to 90%. This interval is the most suitable to perform simulations. The phase transition phenomenon also holds with more than 2 protocols. The more the number of protocols is high, the more the phase transition is shifted to the left. If there are few feasible paths (for small p), the probability that the minimum cost ones involve loops is high. However, this probability quickly decreases. The phase transition phenomenon can be seen in [21]. But the results consider only loopless paths and deals with technologies rather than adaptation functions.

Figure 2b shows the evolution of feasible path existence probability according to p in a scale-free topology. It seems that the probability is independent from the network size, or varying too slowly depending on it. Whether for 50 or 200 nodes, the probability grows from 0.1% to 99% when $p \in [0.01, 0.15]$. The phase transition is shifted to the left compared to the one observed on Topologies $T1$ and $T2$. This can be explained by the fact that the diameter of a scale-free network is asymptotically $\frac{\log n}{\log \log n}$ [29], which, for the chosen values of n , is much smaller than the diameter of Topologies $T1$ and $T2$. And a simple argument shows that the probability of a k nodes path to be feasible is bounded by a value decreasing in k . Consider for simplicity that $\lambda = 2$. There is 12 possible adaptation functions, each one is available on a node with probability p . Thus, the probability that any node is *void*, *i.e.*, without any adaptation function, is $(1-p)^{12}$. The probability for a k nodes path to be feasible is less than $(1 - (1-p)^{12})^k$, which tends to 0. This means that, for fixed p , the longer is a path, the smaller is its feasibility probability.

3) Simulation results on real topologies

Our algorithm is compared to a classical BFS which explores all possible paths until reaching the destination. During

the exploration process, all *dominated*⁴ paths are deleted. BFS can be seen as a version of the algorithm in [11] where the bandwidth constraint is relaxed. The first results showed that BFS algorithm is extremely slow even for small values of p (processing time of the order of several hours). It was impossible to perform a comparison with our algorithm. Due to this tremendous running time, we fixed a maximum length to the explored paths by BFS algorithm. If a path exceeds 10 hops (resp. 14 hops) on Topology $T1$ (resp. $T2$), it is deleted and no more considered. We performed 100 runs for each value of p and averaged the processing time.

Figure 3a shows the processing time of PDA and BFS algorithms on Topologies $T1$ and $T2$ according to the values of p . For small values of p (< 0.22 for $T1$ and < 0.04 for $T2$) BFS algorithm is faster than PDA. However, the processing time of BFS hugely increases. We cannot put it on Figure 3a because it would be unreadable. For example, the processing time of BFS algorithm on Topology $T2$ for $p = 0.24$ is more than 14 minutes, while that of PDA algorithm is 10s. On Topology $T1$, for $p = 0.38$, the processing time of BFS is more than 7 minutes, while that of PDA is 7s.

4) Simulation results on random scale-free topologies

Unlike fixed real topologies, random scale-free topologies allow to study the algorithm efficiency according to the network size, which is the goal of this section. As in the previous case, simulating BFS without path length restriction is impossible due to a prohibitive processing time. Instead, we performed simulations of BFS with maximum path length k from 5 to 9. Paths exceeding this length are deleted and no more considered. Moreover, BFS is only required to find a feasible path, not necessarily the minimum cost one.

Figure 3b shows the processing time of PDA and BFS with different maximum path lengths according to the network size (number of nodes). Figure 3c shows the proportion of feasible paths found by the different BFSs according to the network size. The probability p is fixed to 0.05. All the results are averaged over 1000 runs. The processing time of BFS with $k = 8$ (resp. 9) is tremendous and the algorithm is stopped when the network size reaches 50 (resp 30) nodes. For instance, the processing time of the latter is 177s for 30 nodes. For a number of nodes ≥ 100 , BFSs with $k = 5, 6$, and 7 are faster than PDA. For instance, when the network size is 160, the processing time of PDA is 1261ms while the processing time of BFS $k = 5$ (resp $k = 6$ and $k = 7$) is 1.23ms (resp. 9.78ms and 1005ms). However, as shown in Figure 3c, when the network contains feasible paths, BFS with $k = 5$ finds only 47% of them. This proportion is 83% for $k = 7$. When the network size grows and k is fixed, this proportion tends to 0% since the probability that a feasible path has length $\leq k$ tends to 0. The advantage of PDA is that it is an exact algorithm, it always finds the optimal path when it exists, while the BFS approach lacks both efficiency and correctness since it can miss existing feasible paths.

⁴In this context, a path dominates another one if they have the same extremities and the same protocol stack, and the first path is shorter.

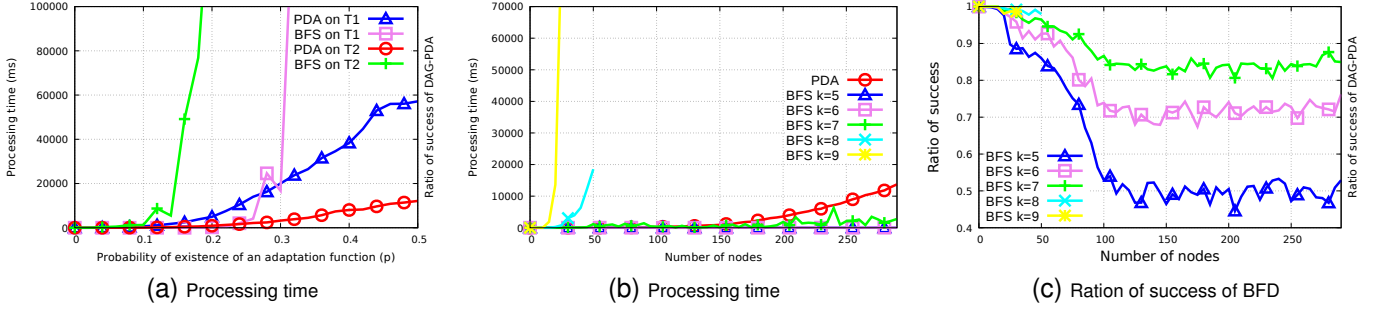


Fig. 3. Processing time of PDA and BFS on scale free topologies and ratio of success of BFS according to the network size.

V. BOUNDS ON FEASIBLE PATH LENGTH

In this section, we distinguish the *cost* of a path, i.e., the sum of the weights of the links and the adaptation functions involved in it, from the *length* of the path, i.e., the number of links involved in it. In both cases, links in loops are counted as many times as they are crossed by the path.

By computing the ℓ -value of the axiom $[Q_0]$, Knuth's algorithm decides if a given multi-layer network contains a feasible path between a specified pair of nodes. If $\ell([Q_0]) = \infty$, then there is no feasible path. Otherwise, there is a feasible path of minimum cost $\ell([Q_0])$. Thus, the existence problem is solvable in $O(\lambda^5 n^2 m)$, which depends only on the network parameters. However, to effectively and explicitly compute the feasible path \mathcal{P} of minimum cost (if any), one needs to generate its trace $\mathcal{T}_{\mathcal{P}}$ and to find the feasible path matching it (Algorithm 3), which are in $O(|\mathcal{T}_{\mathcal{P}}|)$ and $O(|\mathcal{T}_{\mathcal{P}}|nm)$ respectively. Note that $|\mathcal{T}_{\mathcal{P}}|$ is also the length of the feasible path of minimum cost. Thus, the complexity of the whole path computation process is in $O(\max\{\lambda^5 n^2 m, |\mathcal{T}_{\mathcal{P}}|nm\})$, which is polynomial in the network parameters and in the minimum cost path length. However, to accurately estimate this complexity, bounds on the value of $|\mathcal{T}_{\mathcal{P}}|$ are needed.

A. Lower bounds on the shortest feasible path

Obviously, if \mathcal{P} does not involve loops, then $|\mathcal{T}_{\mathcal{P}}|$ is upper bounded⁵ by m . Otherwise, can the length of the minimum cost path be upper bounded by a linear function of n or m ?

Unfortunately, the answer is no. Consider the example depicted on Figure 4. The nodes (except S and D) are partitioned into two sets. For convenience and without loss of generality, we suppose that n is even. Nodes U_1, \dots, U_k where $k = \frac{n}{2} - 1$ form a loop. Each node in this loop is passive (only able to transmit passively any protocol), except the last one which is able to encapsulate any protocol in protocol b . Nodes V_1, \dots, V_k connect U_1 to D via a directed path. All these nodes are able to decapsulate any protocol from protocol b (thus they delete an occurrence of protocol b from the top of the protocol stack). To reach D , a feasible path should cross the sequence of nodes V_1, \dots, V_k and delete an occurrence of b at each of these nodes. Thus, when reaching node V_1 , the

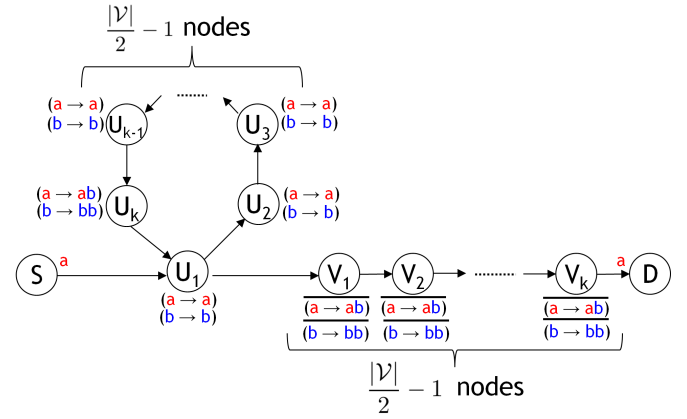


Fig. 4. Example of multi-layer network in which the only feasible path is of length $\Omega(n^2)$.

protocol stack of a feasible path should contain k occurrences of b at its top. However, each crossing of the loop U_1, \dots, U_k adds only one occurrence of b on the top of the stack (at U_k). Therefore, a feasible path should cross the loop k times before reaching V_1 . The length of the shortest (and actually the only one) feasible path is $k^2 + k + 2 = \frac{n^2}{4} - \frac{n}{2} + 2$.

Proposition 2 generalizes this result.

Proposition 2. For any number of nodes n and any set of protocols $\mathcal{A} = \{a_1, \dots, a_\lambda\}$, there is a multi-layer network \mathcal{N} with two specified nodes S and D such that the shortest feasible path from S to D contains $\Omega\left(\left(\frac{n}{\lambda}\right)^\lambda\right)$ links.

Proof. Let \mathcal{N} be a multi-layer network and suppose without loss of generality that $n = \lambda k + 1$ for some integer k . We first construct the underlying graph as following: Create a directed path containing all the nodes where the first one is S and the last one is D . The node $U_{i,j}$ ($0 \leq j < k$) denotes the $ki + j + 1$ th node in the path. Thus D is denoted by $U_{\lambda,0}$. Then create a link from node $U_{i,0}$ to S for all $i \leq k - 1$. The adaptation functions are allocated as follows:

- Node S emits protocol a_1 and nodes $U_{0,1}, \dots, U_{0,k-1}$ are passive and can only transmit any protocol as received,
- Nodes $U_{i,1}, \dots, U_{i,k-1}$ ($1 \leq i \leq \lambda - 1$) are only able to decapsulate protocol a_{i+1} from any protocol,
- Nodes $U_{i,0}$ ($1 \leq i \leq \lambda - 1$) are only able to encapsulate any protocol in protocol a_{i+1} .

⁵Since the weights and other possible QoS parameters depend on the links (and the adaptation functions) but not on vertices, paths involving several times the same vertices are not an issue and are considered as loopless.

A path must have the sequence a_{i+1}^{k-1} at the top of the protocol stack⁶ to cross the nodes $U_{i,1}, \dots, U_{i,k-1}$ ($1 \leq i \leq \lambda - 1$) and reach $U_{i+1,0}$. To do so, it must cross the node $U_{i,0}$ $k - 1$ times because it is the only node able to add an occurrence of a_{i+1} on the top of the protocol stack. And at each time it reaches $U_{i,0}$ (except the last one), it must cross the link $(U_{i,0}, S)$ again. Let ϕ_i be the shortest path length (in links) to reach $U_{i,0}$. Then reaching $U_{i,0}$ $k - 1$ times requires $(k - 2)(\phi_i + 1) + \phi_i$ links. Thus, reaching $U_{i+1,0}$ requires $(k - 2)(\phi_i + 1) + \phi_i + k$ links. Rearranged, this gives the recurrent equation

$$\phi_{i+1} = (k - 1)\phi_i + 2k - 2 \text{ with } \phi_1 = k \quad (1)$$

Resolving the equation gives $\phi_\lambda = \Omega(k^\lambda)$. Replacing k by $\frac{n-1}{\lambda}$ concludes the proof.

The previous result gives an example of a shortest feasible path that is exponential in the number of protocols but polynomial in the number of nodes. However, it is possible to obtain a superpolynomial shortest path with only two protocols. It is well known that the shortest word generated by a CFG can be of exponential length in the CFG size. On the other hand, there is a polynomial relation between the size of a (W)CFG and an equivalent (W)PDA. This suggests that the shortest feasible path can be exponential. However, the classical algorithm that converts a CFG into a PDA produces a PDA with a fixed number of states (generally only one - see [24] for more details). Thus, the resulting (W)PDA is unlikely to be converted into a multi-layer network. Instead, we build a multi-layer network where the shortest path length is superpolynomial. Proposition 3 details this construction.

Proposition 3. *For any number of nodes n , there is a multi-layer network having 2 protocols with two specified nodes S and D , such that the shortest feasible path from S to D contains $\Omega(\sqrt{n}2^{\Omega(\sqrt{n})})$ links.*

Proof. Let \mathcal{N} be the multi-layer network built in the proof of Proposition 2 with λ protocols. We will convert it into a multi-layer network \mathcal{N}' with two protocols a and b . The key idea of the proof is to emulate a protocol a_{i+1} with the sequence of protocols $ab^i a$. To do so, each node $U_{i,0}$ ($1 \leq i \leq \lambda - 1$) is replaced by a sequence of $i + 2$ nodes $U_{i,0}^0 \dots U_{i,0}^{i+1}$ where $U_{i,0}^0$ and $U_{i,0}^{i+1}$ are able to encapsulate any protocol in protocol a and the nodes $U_{i,0}^1 \dots U_{i,0}^i$ are able to encapsulate any protocol in protocol b . Thus, crossing this sequence push $ab^i a$ on the top of the protocol stack. Analogously, each node $U_{i,j}$ ($1 \leq i \leq \lambda - 1$ and $1 \leq j \leq k - 1$) is replaced by a sequence of $i + 2$ nodes $U_{i,j}^0 \dots U_{i,j}^{i+1}$ where $U_{i,j}^0$ and $U_{i,j}^{i+1}$ are able to decapsulate protocol a from any protocol and the nodes $U_{i,j}^1 \dots U_{i,j}^i$ are able to decapsulate protocol b from any protocol. To cross this sequence of nodes, the top of the stack should contain $ab^i a$. Now, let ϕ_i be the shortest path length (in links) to reach $U_{i,0}^{i+1}$. The recurrent relation (1) becomes

$$\phi_{i+1} = (k - 1)\phi_i + (i + 3)k - 2 \text{ with } \phi_1 = k + 2$$

⁶The notation a^k stands for $\underbrace{aa \dots a}_{k \text{ times}}$.

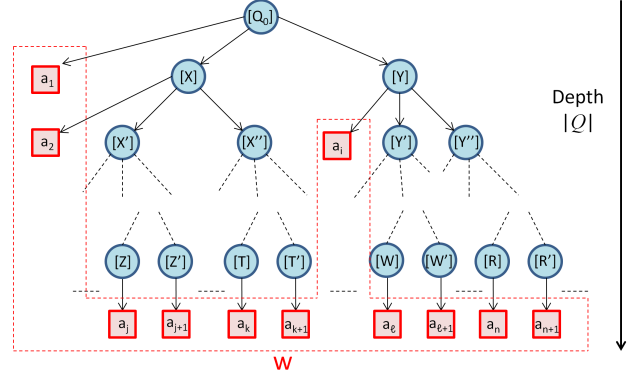


Fig. 5. Derivation tree of a word corresponding to a feasible path.

Resolving the equation gives $\phi_\lambda = \Omega(\lambda k^\lambda)$. Replacing λ by $\frac{n-1}{k}$ and setting $k = 3$ gives $\Omega(n3^{\frac{n}{3}})$. However, here n is the number of nodes of \mathcal{N} while the path we exhibit is in \mathcal{N}' . The number of nodes in \mathcal{N}' is $n' = O(n^2)$, thus $n = \Omega(\sqrt{n'})$. The shortest feasible path length in \mathcal{N}' is in $\Omega(\sqrt{n'}3^{\Omega(\sqrt{n'})}) = \Omega(\sqrt{n'}2^{\Omega(\sqrt{n'})})$ \square

B. Upper bound on the shortest feasible path length

Proposition 4. *For any multi-layer network \mathcal{N} , the shortest feasible path (if any) between two nodes contains at most $O(2^\psi)$ hops, where $\psi = (\lambda + 1)\lambda^2 n^2$.*

Proof. We prove that the shortest word (if any) generated by a CFG corresponding to any multi-layer network contains at most $2^{|\mathcal{Q}|} - 1$ terminals, where $|\mathcal{Q}|$ is the number of nonterminals in the CFG. The proof is similar to the one showing bounds on the length of the shortest word generated by a CFG in Chomsky normal form [24].

Let \mathcal{N} be any multi-layer network and let (S, D) be the source and the destination of a feasible path in \mathcal{N} . Recall that the CFG corresponding to \mathcal{N} and produced by Algorithm 2 contains production rules of the following forms:

- 1) $[X] \rightarrow a$
- 2) $[X] \rightarrow a[Y]$
- 3) $[X] \rightarrow a[Y][Z]$

Where a is a terminal or the empty word ϵ , and $[X], [Y], [Z]$ are nonterminals. Let w be the trace of the shortest feasible path between S and D in \mathcal{N} . The word w is then the shortest one generated by the CFG. The derivation tree of w is a tree where the axiom is the root, the internal nodes are nonterminals and the leafs are terminals or the empty word. A nonterminal $[X]$ is the parent of a nonterminal $[Y]$ if the derivation of w contains a production rule where $[X]$ is the left side and $[Y]$ appears in the right side. $[X]$ is the parent of a leaf a in the same way. For example, if the production rule $[X] \rightarrow a[Y][Z]$ is used to derive w , then $[X]$ will be the parent of the leaf a and the nodes $[Y]$ and $[Z]$. Such a tree is a derivation tree of w means that the sequence of leafs of the tree from left to right corresponds exactly to w .

To maximize the size of the derivation tree of w , suppose that it contains only production rules of form 3, except those of the bottom of the tree, which must be in form 1 (otherwise w

would contain nonterminals). Furthermore, suppose that in all the production rules of the derivation tree, a is never the empty word. Figure 5 illustrates the shape of such a tree and the way the word w is derived. This is the biggest possible derivation tree for w . The key argument is the following: Each branch in the derivation tree contains each nonterminal of the CFG at most once. If a nonterminal appears twice or more in a branch, one could generate a shorter word by using this nonterminal only once, and this would contradict the fact that w is the shortest word. Thus, the biggest possible derivation tree for w is a complete tree where each internal node is parent of one leaf (terminal) and two internal nodes (nonterminals), except the last node of each branch which is parent of only a leaf. The depth of the tree is at most $|Q|$, and the number of internal nodes is at most $2^{|Q|} - 1$. Since each internal node gives at most one terminal, then the length of w is at most $2^{|Q|} - 1$. The number of nonterminals in the CFG is at most $|\Gamma| \times |\mathcal{S}|^2 = (\lambda + 1)\lambda^2 n^2$ (see Algorithm 2 and its complexity). Thus, the path length is bounded by $2^{(\lambda+1)\lambda^2 n^2} - 1$. \square

VI. ADDRESSING BANDWIDTH CONSTRAINT

This section studies the complexity of path computation under bandwidth constraint and proposes heuristic solutions to resolve the problem.

A. Problem formalization

For Traffic Engineering purposes, a feasible path may be constrained by a minimal bandwidth. But it is possible that feasible paths in a multi-layer network involve loops (i.e., involving the same link several times but using different protocols). It implies that the bandwidth constraint is no longer prunable: Even if the links with not enough bandwidth are deleted by topology filtering prior to path computation, other links can have enough bandwidth if they are selected once but not if more. For example, if a link has a capacity of 10Gbps and the bandwidth constraint is 5Gbps, then this link cannot be crossed more than twice. The (optimization) problem of computing the minimum cost path in a multi-layer network under bandwidth constraint is defined as follows:

$$\begin{aligned} \min h(\mathcal{P}) &= \sum_{(U,f,V) \in \mathcal{P}} h(U, f, V) \\ \text{s.t. } &\left\{ \begin{array}{l} \mathcal{P} \text{ is a feasible path between } S \text{ and } D \\ \min_{E \in \mathcal{P}} \frac{q_b(E)}{nb(E)} \geq q_b^{\min} \end{array} \right. \quad (2) \end{aligned}$$

where $nb(E)$ is the number of times a link E is crossed by path \mathcal{P} , $q_b(E)$ is the bandwidth capacity of E and q_b^{\min} is the bandwidth constraint.

B. Path computation complexity under bandwidth constraint

The bandwidth constraint impacts the complexity of feasible path computation. In a single-layer network, computing a path under bandwidth constraint is trivial: It suffices to prune all the links without enough bandwidth. This is no longer possible in a multi-layer network. In fact, the decision problem is NP-hard as shown by Kuipers and (Free)k Dijkstra [11]. But this

proof does not work on symmetric directed graphs⁷. However, most communication networks are symmetric. We show that the decision version of the problem remains NP-hard even with two protocols and in a symmetric graph. Consider the following problem:

Problem (2'). Given a multi-layer network $\mathcal{N} = (\mathcal{G} = (\mathcal{V}, \mathcal{E}), \mathcal{A}, \mathcal{F}, h)$, a function assigning to each link $E \in \mathcal{E}$ an available bandwidth $q_b(E)$, a bandwidth constraint q_b^{\min} and a pair S and D of nodes in \mathcal{V} . Is there a feasible path from S to D satisfying the bandwidth constraint?

Proposition 5. *Problem (2') is NP-hard with two protocols even if $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ is a symmetric directed graph.*

Proof. First consider the problem of finding a Hamiltonian path in a symmetric directed graph between two nodes S' and D' . Call this problem SYM-HAM. SYM-HAM is NP-complete (a polynomial reduction from the classical Hamiltonian path problem to SYM-HAM is trivial).

Now we provide a polynomial reduction from SYM-HAM to Problem (2') restricted to symmetric directed graph and two protocols. Given an instance of SYM-HAM, i.e., a symmetric directed graph $\mathcal{H} = (\mathcal{V}', \mathcal{E}')$ and a pair of nodes (S', D') , we build an instance of Problem (2'), i.e., a network $\mathcal{N} = (\mathcal{G}, \mathcal{A}, \mathcal{F}, h)$ and a pair of nodes (S, D) as following:

Step 1: Splitting the nodes. For each node $U' \in \mathcal{V}'$, four nodes U_1, U_2, U_3 and U_4 are created in \mathcal{G} . Links (U_i, U_{i+1}) and (U_{i+1}, U_i) are created for $i = 1, \dots, 3$. For each link $(U', V') \in \mathcal{E}'$, a link (U_1, V_1) is created in \mathcal{G} . This step is illustrated on Figure 6.

Step 2: Adding a tail. $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ is augmented by a set $\mathcal{C} = \{C_0, \dots, C_{n+1}\}$ of nodes, where $C_0 = S$ is the source node. There are a link (C_i, C_{i+1}) and a link (C_{i+1}, C_i) for $i = 0, \dots, n$. Moreover, there is also a link from C_{n+1} to S_1 (the first node resulting from the splitting of S') and conversely. Figure 7 shows this construction. Finally, two nodes X and D are added, as well as the four links (D_1, X) , (X, D_1) , (X, D) and (D, X) (recall that D_1 is the first node resulting from the splitting of D' , see Step 1).

Step 3: Allocating the adaptation functions and available bandwidth. All the links have available bandwidth 1. The bandwidth constraint is set to 1. Thus, any feasible path must cross a link at most once. There is no possible loop. Let the set of protocols be $\mathcal{A} = \{a, b\}$. Node S emits packets of protocol a . For $i = 1 \dots, n$, each node C_i in the tail can encapsulate protocol a in itself. Node C_{n+1} can only encapsulate a in b . For each node $U' \in \mathcal{V}'$, node U_1 can encapsulate any protocol in b . Node U_2 can either decapsulate protocol b from itself or passively transmit protocol a . Node U_3 can either decapsulate protocol a from b or passively transmit protocol a . Node U_4 is able to decapsulate protocol a from itself. Finally, node X can decapsulate protocol a from b . Table I summarizes the allocation of the adaptation functions.

Now, we prove that there is a Hamiltonian path from S' to D' in \mathcal{H} if and only if there is a feasible path from S to D in \mathcal{N} that satisfies the bandwidth constraint. First, assuming that

⁷We define a symmetric directed graph as a graph where a link (U, V) exists if and only if the reverse link (V, U) exists.

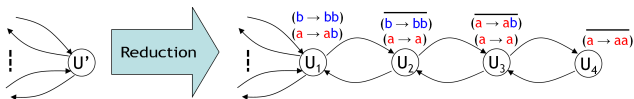


Fig. 6. Reduction from SYM-HAM to feasible path under bandwidth constraint (node splitting).

Node	Adaptation functions
$C_i, i = 1 \dots n$	$(a \rightarrow aa)$
C_{n+1}	$(a \rightarrow ab)$
U_1 s.t. $U' \in \mathcal{V}'$	$(b \rightarrow bb), (a \rightarrow ab)$
U_2 s.t. $U' \in \mathcal{V}'$	$(b \rightarrow bb), (a \rightarrow a)$
U_3 s.t. $U' \in \mathcal{V}'$	$(a \rightarrow ab), (a \rightarrow a)$
U_4 s.t. $U' \in \mathcal{V}'$	$(a \rightarrow aa)$
X	$(a \rightarrow ab)$

TABLE I

THE ADAPTATION FUNCTIONS AVAILABLE ON THE NODES IN THE POLYNOMIAL REDUCTION.

there is a Hamiltonian path from S' to D' in \mathcal{H} , we construct a feasible path \mathcal{P} in \mathcal{N} as follows: Starting from S in \mathcal{N} , \mathcal{P} crosses the tail and each C_i ($i = 1 \dots n$) adds an occurrence of protocol a in the stack of encapsulated protocols. Then crossing C_{n+1} adds b as current protocol. Thus, at the end of the tail, there are $n + 1$ encapsulated protocols a (the one emitted by S and n occurrences added in the tail) and the current protocol is b . Following the same node order as in the Hamiltonian path, replace each occurrence of a node $U' \in \mathcal{V}'$ (including S' and D') in the Hamiltonian path by the sequence:

$$\begin{aligned}
 &U_1(b \rightarrow bb)U_2(\overline{b \rightarrow bb})U_3(\overline{a \rightarrow ab})U_4(\overline{a \rightarrow aa})U_3(a \rightarrow a) \\
 &U_2(a \rightarrow a)U_1(a \rightarrow ab)
 \end{aligned} \tag{3}$$

Thus, at node U_1 an encapsulation of protocol b occurs, at U_2 protocol b is decapsulated, at U_3 it is decapsulated again, and at U_4 protocol a is decapsulated. Path \mathcal{P} then crosses passively nodes U_3 and U_2 , and finally encapsulates protocol b at U_1 . Thus, at each time the path crosses an occurrence of Sequence (3), then one occurrence of protocol a is removed from the protocol stack. Crossing all U_4 s.t. $U' \in \mathcal{H}$ removes all encapsulated occurrences of protocol a except the first one. When the path leaves D_1 to reach node X , the current protocol is b and there is a last occurrence of protocol a which is encapsulated. Finally, node X decapsulates protocol a from protocol b and node D receives protocol a as emitted by S . Thus, \mathcal{P} is a feasible path, and each link is crossed at most once, the bandwidth constraint is satisfied.

Conversely, we show that from any feasible path \mathcal{P} satisfying the bandwidth constraint in \mathcal{N} , one can extract a Hamiltonian path between S' and D' in \mathcal{H} . A feasible path must cross all nodes U_4 s.t. $U' \in \mathcal{V}'$ in order to decapsulate all occurrences of protocol a encapsulated when crossing the tail. Thus, it involves Sequence (3) for all $U' \in \mathcal{V}'$. By removing the tail part and the nodes X and D from \mathcal{P} and replacing each occurrence of Sequence (3) by the corresponding node U' , the

resulting path starts from S' and crosses all the nodes in \mathcal{H} before reaching D' . The only problem is the possibility that there are other sequences than Sequence (3) in the remaining path. There are two possible cases:

- An *incomplete* Sequence (3) where U_4 is not reached (e.g., $U_1fU_2f'U_3f''U_2f'''U_1$): This cannot happen because such a sequence forbids to reach U_4 later, and thus one encapsulated occurrence of protocol a is never decapsulated and \mathcal{P} cannot be feasible. Such a sequence cannot occur after an occurrence of Sequence (3) on the same nodes because if a node U_i ($i = 2, 3$) is reached in a Sequence (3) it cannot be reached again due to the bandwidth constraint.
- A sequence $U_1fV_1f'W_1$: Let \mathcal{P} be a feasible path from S to D containing a sequence $U_1fV_1f'W_1$ (where U_1 and W_1 may be the same node). These three nodes can only encapsulate protocol a or b in protocol b . Thus, after crossing such a sequence, there are three occurrences of protocol b on the top of the protocol stack. However, in network \mathcal{N} , there is no possible sequence of nodes and adaptation functions able to decapsulate protocol b three consecutive times. Thus, \mathcal{P} is not feasible.

Thus, if a feasible path exists, then it contains only one occurrence of Sequence (3) for each node $U' \in \mathcal{V}'$. Replacing each Sequence (3) by the corresponding node in \mathcal{V}' induces a Hamiltonian path in \mathcal{H} . This concludes the proof. \square

Remark. In the conference version of this paper [13], Proposition 1 incorrectly states that Problem (2') is NP-complete, and thus being in NP. This claim was based on the supposition that the shortest feasible path length should be polynomial in the size of the network. However, in the light of the results stated in Section V, this claim is false and the shortest path cannot be a polynomial certificate.

Corollary 1. *Problem (2) is not approximable.*

Corollary 2. *Problems (2) and (2') are polynomial if the graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ is a Directed Acyclic Graph (DAG).*

C. DAG Heuristic

As seen in Section IV-B2, minimum cost feasible paths involving loops are infrequent (for $p > 20\%$). Combining this fact with Corollary 2 suggests a heuristic to compute feasible path under bandwidth constraint: Convert the network into a DAG and perform the PDA algorithm to compute a minimum cost feasible path.

DAG Conversion. The network is converted into a DAG in the following way:

- 1) Set the number 0 to node S and $n - 1$ to node D (recall that S and D are the extremities of the graph diameter);
- 2) Perform a BFS algorithm starting from node S and number the nodes in the visit order. The nodes at the same distance from S are visited randomly, thus performing several times this heuristic does not always give the same node numbering and the same DAG;
- 3) Delete all the links that start at a node and end at a node with a smaller number.

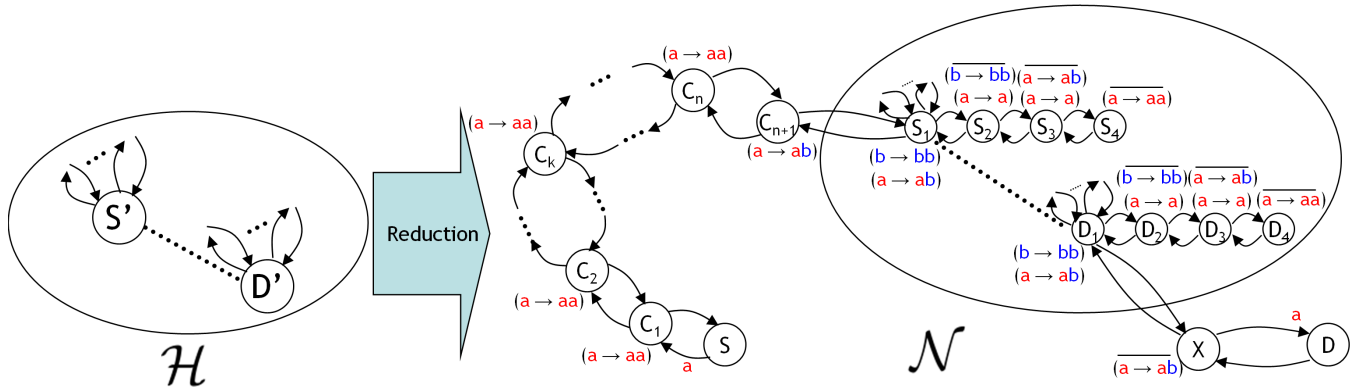


Fig. 7. Reduction from SYM-HAM to feasible path under bandwidth constraint (graph transformation).

The DAG heuristic is as follows:

- 1) Convert the network into a DAG;
- 2) Prune the links without enough bandwidth;
- 3) Perform the PDA algorithm of Section IV.

D. Simulations

We study the efficiency of the DAG heuristic (called DAG-PDA) and compare it with the algorithm of Kuipers and Dijkstra [11]. The latter is an exact (and thus exponential) algorithm that performs a BFS and explores all the paths that are not dominated and that satisfy the bandwidth constraint. As in Section IV-B3, the BFS algorithm is slow. Thus, we also compare our algorithm with DAG-BFS algorithm, where the network is converted into a DAG before performing the BFS. The simulation conditions (parameters, topology, number of runs, etc.) are the same as in Section IV-B3. The bandwidth capacity of the links is randomly and uniformly selected in the set $\{1, 2, \dots, 10\}$. The bandwidth constraint is set to 2.

1) Comparison of the feasibility ratio

Converting the network topology into a DAG deletes some feasible paths in the original network. We measure how much feasible paths are lost by comparing the probability of feasible path existence before and after the DAG conversion according to the probability of existence of adaptation functions (p).

Figure 8a shows that the probability of feasible path existence is shifted to the right after the DAG conversion. The ratio $\frac{\text{Probability of feasible path existence in DAG } T_i}{\text{Probability of feasible path existence in } T_i}$ ($i = 1, 2$) is clearly decreasing and is less than 50% if $p > 0.34$, which is important but balanced by the improvement of the processing time.

2) Comparison of the processing time

Figure 8b shows the processing time of DAG-PDA, DAG-BFS and BFS algorithms on both topologies according to the probability of existence of an adaptation function. BFS algorithm is slow even for small values of p . For $p < 0.3$ (resp. 0.4) on Topology $T1$ (resp. $T2$), DAG-BFS is faster than DAG-PDA. Beyond these values, the processing time of DAG-BFS hugely increases. For example, for $p = 0.5$, the processing time of DAG-BFS is more than 35 minutes on Topology $T1$ and more than 53 minutes on Topology $T2$, while that of DAG-PDA is 3.8 seconds on $T1$ and 24 seconds

on $T2$. These results show that the DAG-PDA algorithm is clearly faster when there is a significant number of adaptation functions, but the exponential DAG-BFS algorithm is faster if there are few of them (for small values of p).

VII. PATH COMPUTATION UNDER QoS CONSTRAINTS

A. Multi-constrained feasible path

Let \mathcal{N} be a multi-layer network. Each link $E = (U, V)$ is associated to a set of k additive QoS metrics $q(E) = (q_1(E), \dots, q_k(E))$ in addition to its available bandwidth $q_b(E)$. These additive metrics can be the delay, logarithm of the packet-loss, etc. Let q_b^{\min} be the bandwidth constraint and $q^{\max} = (q_1^{\max}, q_2^{\max}, \dots, q_m^{\max})$ be a vector of QoS constraints, the problem of computing a minimum cost feasible path under these constraints is formalized as:

$$\min h(\mathcal{P}) = \sum_{(U, f, V) \in \mathcal{P}} h(U, f, V)$$

$$s.t. \begin{cases} \mathcal{P} \text{ is a feasible path between } S \text{ and } D \\ \min_{E \in \mathcal{P}} \frac{q_b(E)}{nb(E)} \geq q_b^{\min} \\ \sum_{E \in \mathcal{P}} (q_i(E) \times nb(E)) \leq q_i^{\max}, i = 1 \dots k \end{cases} \quad (4)$$

The problem of QoS multi-constrained path computation (on a single layer) is well studied and it is well-known that the decision version is NP-complete, even with 2 additive and/or multiplicative constraints [30]. It is a particular case of Problem 4, corresponding to the case where there is only one protocol and passive transitions. Thus the decision version associated to Problem 4 is NP-hard.

B. ML-SAMCRA

Van Mieghem and Kuipers [12] proposed the Self-Adaptive Multiple Constraints Routing Algorithm (SAMCRA) to resolve the problem of multi-constrained path computation on a single layer. It computes paths under several (additive) QoS constraints but it ignores the feasibility constraint as defined in our paper. The algorithm has exponential time worst case complexity but the authors showed that it exhibits a reasonable processing time in practice. We propose to adapt SAMCRA to the multi-layer context.

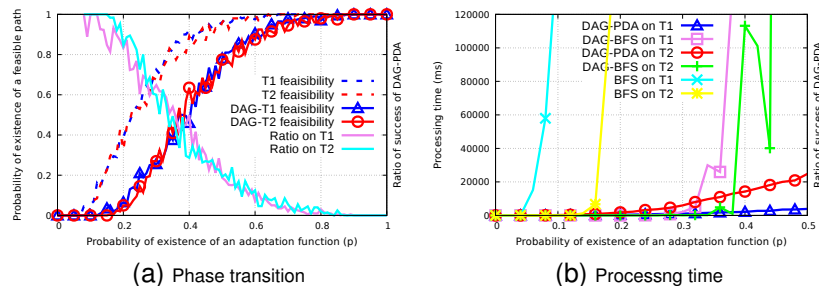


Fig. 8. Phase transition (before and after DAG conversion) and comparison of the processing time of DAG-PDA, DAG-BFS and BFS.

1) The main concepts of SAMCRA

The idea of SAMCRA is to maintain a path list from the source node S to all other nodes until reaching the destination node D . It progressively removes the paths that do not comply with the QoS constraints. The main concepts of SAMCRA are:

- *Non-linear path length*: It reduces the solution space to scan but the algorithm can apply with any metric. Hence, it is not a strict requirement.
- *The k -shortest path algorithm*: The k -shortest path algorithm maintains the list of the paths that are not (yet) removed from the path list.
- *Non-dominance*: A multi-constrained path \mathcal{P} dominates another path \mathcal{P}' if $\forall i, \sum_{E \in \mathcal{P}} q_i(E) \leq \sum_{E \in \mathcal{P}'} q_i(E)$ (i.e., if \mathcal{P} is better than \mathcal{P}' for each QoS parameter). A path \mathcal{P} is non-dominated if there is no path which dominates it. The concept of non-dominance induces a partial order over the paths. It avoids the exploration of several paths thus substantially reducing the average complexity of SAMCRA.

The only concept that is impacted by the multi-layer context is the non-dominance. It must be redefined to meet the path feasibility constraint and to take into account possible loops.

2) Extension of the non-dominance definition

A multi-layer path is characterized by its nodes but also by its protocol stack at the end node. Thus in the algorithm path list, each path should be stored with its protocol stack at its final node. A path can involve the same link several times. Before checking if it complies with some QoS parameters, the parameters of each link should be multiplied by the number of times this link is involved in the path. The bandwidth constraint is not prunable in multi-layer context, the new non-dominance definition should take it into account.

A path \mathcal{P} dominates a path \mathcal{P}' if and only if:

- $\min_{E \in \mathcal{P}} \frac{q_b(E)}{nb_{\mathcal{P}}(E)} \geq \min_{E \in \mathcal{P}'} \frac{q_b(E)}{nb_{\mathcal{P}'}(E)}$
- $\sum_{E \in \mathcal{P}} q_i(E) \times nb_{\mathcal{P}}(E) \leq \sum_{E \in \mathcal{P}'} q_i(E) \times nb_{\mathcal{P}'}(E)$
 $\forall i = 1, \dots, k$
- \mathcal{P} and \mathcal{P}' have the same final node;
- \mathcal{P} and \mathcal{P}' have the same protocol stack at this node.

Where $nb_{\mathcal{P}}(E)$ (resp. $nb_{\mathcal{P}'}(E)$) is the number of times the link E is involved in path \mathcal{P} (resp. \mathcal{P}'). According to this new definition of non-dominance, ML-SAMCRA explores all the possible paths until reaching the destination node with satisfactory QoS parameters. Along the exploration, it removes all paths that are dominated or not feasible.

C. Some simulation results

We studied the efficiency of ML-SAMCRA through simulations and checked if it is as scalable in a multi-layer context as SAMCRA is in a single layer context. The same simulation methodology as in Section IV-B is applied. The results showed that for $p > 0.08$ (resp. 0.10) on Topology $T1$ (resp. $T2$) the processing time hugely increases (more than 1 minutes). Clearly, ML-SAMCRA does not scale above these values. The main possible reason is that the paths are less comparable in term of the new non-dominance definition: They should have the same protocol stack. As there are less dominated paths, the algorithm complexity increases. The second possible reason is that taking into account loops increases the number and the length of the paths, which in turn increases the algorithm complexity. Thus, path computation under QoS constraints in multi-layer networks is more complex than in single layer networks. Exact algorithms are suitable only for small instances.

VIII. CONCLUSION

Most of carrier-grade networks manage their different layers thanks to separate control planes. Designing a unified control plane would allow the network resources to be optimized and the operational management costs to be reduced. One key problem to address is path computation taking into account the protocol heterogeneity and the multi-layer context dealing with encapsulation, conversion and decapsulation of protocols. This paper tackles this issue by partitioning it into three cases: Path computation without bandwidth constraint, under bandwidth constraint and under additive QoS constraints. For the first case, we widely generalized polynomial algorithms in the state of the art and decreased their complexity. Through simulations, we showed that they are faster than the previous approaches in the literature. We also provided lower and upper bounds on the shortest path length. For the second case, we obtained several time complexity results and proposed efficient heuristics. Finally, we designed the first algorithm to resolve the third case. In future works, we plan to design heuristics to deal with additive QoS metrics, as the exact approach is not scalable. The problem of efficient generation of random topologies being widely open, it would be interesting to analytically study the phase transition phenomenon in order to generate topologies having a suitable number of feasible paths.

REFERENCES

- [1] L. Martini, E. Rosen, N. El-Aawar, and G. Heron, "RFC4448 - Encapsulation Methods for Transport of Ethernet over MPLS Networks," 2008.
- [2] F. Baker, X. Li, C. Bao, and K. Yin, "RFC6144 - Framework for IPv4/IPv6 Translation," 2011.
- [3] S. Bryant and P. Pate, "RFC3985 - Pseudo Wire Emulation Edge-to-Edge (PWE3) Architecture," 2005.
- [4] H. Song, "Protocol-oblivious forwarding: Unleash the power of sdn through a future-proof forwarding plane," in *Proceedings of the second ACM SIGCOMM workshop on Hot topics in software defined networking*. ACM, 2013, pp. 127–132.
- [5] P. Bosshart, D. Daly, G. Gibb, M. Izzard, N. McKeown, J. Rexford, C. Schlesinger, D. Talayco, A. Vahdat, G. Varghese *et al.*, "P4: Programming protocol-independent packet processors," *ACM SIGCOMM Computer Communication Review*, vol. 44, no. 3, pp. 87–95, 2014.
- [6] S. Das, G. Parulkar, N. McKeown, P. Singh, D. Getachew, and L. Ong, "Packet and circuit network convergence with openflow," in *Optical Fiber Communication Conference*. Optical Society of America, 2010.
- [7] L. Liu, D. Zhang, T. Tsuritani, R. Vilalta, R. Casellas, L. Hong, I. Morita, H. Guo, J. Wu, R. Martínez *et al.*, "Field trial of an openflow-based unified control plane for multilayer multigranularity optical switching networks," *Lightwave Technology, Journal of*, vol. 31, no. 4, pp. 506–514, 2013.
- [8] S. Agarwal, M. S. Kodialam, and T. V. Lakshman, "Traffic engineering in software defined networks," in *Proceedings of the IEEE INFOCOM 2013, Turin, Italy, April 14-19, 2013*, 2013, pp. 2211–2219.
- [9] S. Li, Y. Shao, S. Ma, N. Xue, S. Li, D. Hu, and Z. Zhu, "Flexible traffic engineering: When openflow meets multi-protocol ip-forwarding," *IEEE Communications Letters*, vol. 18, no. 10, pp. 1699–1702, 2014.
- [10] M. L. Lamali, H. Pouyllau, and D. Barth, "Path computation in multi-layer multi-domain networks: A language theoretic approach," *Computer Communications*, vol. 36, no. 5, pp. 589–599, 2013.
- [11] F. A. Kuipers and F. Dijkstra, "Path selection in multi-layer networks," *Computer Communications*, 2009.
- [12] P. V. Mieghem and F. A. Kuipers, "Concepts of exact QoS routing algorithms," *IEEE/ACM Trans. Netw.*, vol. 12, no. 5, pp. 851–864, 2004.
- [13] M. L. Lamali, N. Fergani, J. Cohen, and H. Pouyllau, "Path computation in multi-layer networks: Complexity and algorithms," in *INFOCOM 2016, San Francisco, CA, USA, April 10-14, 2016*, 2016, pp. 1–9.
- [14] A. Farrel, J. Vasseur, and J. Ash, "RFC4655 - A Path Computation Element (PCE)-Based Architecture," 2006.
- [15] M. Bocci and S. Bryant, "RFC5659 - An Architecture for Multi-Segment Pseudowire Emulation Edge-to-Edge," 2009.
- [16] F. Dijkstra, J. V. der Ham, P. Grosso, and C. de Laat, "A path finding implementation for multi-layer networks," *Future Generation Comp. Syst.*, vol. 25, no. 2, pp. 142–146, 2009.
- [17] I. Chlamtac, A. Faragó, and T. Zhang, "Lightpath (Wavelength) Routing in Large WDM Networks," *IEEE Journal on Selected Areas in Communications*, vol. 14, no. 5, pp. 909–913, 1996.
- [18] H. Zhu, H. Zang, K. Zhu, and B. Mukherjee, "A novel generic graph model for traffic grooming in heterogeneous WDM mesh networks," *IEEE/ACM Trans. Netw.*, vol. 11, no. 2, pp. 285–299, 2003.
- [19] S. Gong and B. Jabbari, "Optimal and Efficient End-to-End Path Computation in Multi-Layer Networks," in *ICC*, 2008, pp. 5767–5771.
- [20] F. Dijkstra, B. Andree, K. Koymans, J. van der Ham, P. Grosso, and C. de Laat, "A multi-layer network model based on ITU-T G.805," *Comput. Netw.*, 2008.
- [21] F. Iqbal, J. van der Ham, and F. Kuipers, "Technology-aware multi-domain multi-layer routing," *Computer Communications*, vol. 62, pp. 85–96, 2015.
- [22] D. E. Knuth, "A Generalization of Dijkstra's Algorithm," *Inf. Process. Lett.*, vol. 6, no. 1, pp. 1–5, 1977.
- [23] I. Petre and A. Salomaa, "Algebraic Systems and Pushdown Automata," in *Handbook of Weighted Automata*, M. Droste, W. Kuich, and H. Vogler, Eds. Springer Publishing Company, Incorporated., 2009, ch. 7, pp. 257–290.
- [24] J. E. Hopcroft, R. Motwani, and J. D. Ullman, "Introduction to automata theory, languages, and computation." Addison-Wesley, 2006.
- [25] M. L. Fredman and R. E. Tarjan, "Fibonacci Heaps and Their Uses in Improved Network Optimization Algorithms," *J. ACM*, vol. 34, no. 3, pp. 596–615, Jul. 1987.
- [26] R. Summerhill, "The new internet2 network," in *6th GLIF Meeting*, 2006, available at <http://www.internet2.edu/products-services/advanced-networking>.
- [27] R. Mahajan, N. Spring, D. Wetherall, and T. Anderson, "Inferring link weights using end-to-end measurements," in *Proceedings of the 2nd ACM SIGCOMM Workshop on Internet measurement*. ACM, 2002, pp. 231–236. [Online]. Available: <http://research.cs.washington.edu/networking/rocketfuel/>
- [28] A.-L. Barabási and R. Albert, "Emergence of scaling in random networks," *science*, vol. 286, no. 5439, pp. 509–512, 1999.
- [29] B. Bollobás and O. Riordan, "The diameter of a scale-free random graph," *Combinatorica*, vol. 24, no. 1, pp. 5–34, 2004.
- [30] Z. Wang and J. Crowcroft, "Quality-of-Service Routing for Supporting Multimedia Applications," *IEEE Journal on Selected Areas in Communications*, vol. 14, no. 7, pp. 1228–1234, 1996.

ACKNOWLEDGMENT

The authors would like to thank Hélia Pouyllau for her participation to the previous version of this work.



Mohamed Lamine Lamali received the M.Sc. degree and Ph.D. degree from University of Versailles, France, in 2010 and 2014, respectively. He was member of Bell Labs France from 2010 to 2016, in the "Advanced Internet Research" department then in the "Mathematics for Dynamic Networks" department. Since sept. 2016, he is associate professor at University of Bordeaux. His work mainly focuses on algorithmic and complexity aspects of networking, applied to Quality of Service routing and path computation in multi-layer networks.



Nasreddine Fergani received a Master's degree in Operations Research from the University of Science and Technology of Algiers, and an Engineering degree in Computer Science from Polytech Tours, France in 2015. He has held an intern position at Bell Labs Paris, a Research assistant position at Queen Mary University of London. He is now part of the Engineering Development Group at MathWorks, Cambridge.



Johanne Cohen has worked at Laboratory LRI, University ParisSud, since 2013. She obtained her PhD degree in Computer Science from the University Paris-Sud in 1998 and her HDR (Habilitation Diriger les Recherches) in 2009. She has been a full time researcher at CNRS since 2001. Her research interests include algorithmics for telecommunication, graph theory, algorithmic game theory, distributed algorithm, complexity and approximation theory. She leads the Graphs, Algorithmic and Combinatorics the Algorithmic team.