Enhancing Resource Utilization of Non-terrestrial Networks Using Temporal Graph-based Deterministic Routing

Keyi Shi, Jingchao Wang, Hongyan Li, Member, IEEE, and Kan Wang

Abstract—Deterministic routing has emerged as a promising technology for future non-terrestrial networks (NTNs), offering the potential to enhance service performance and optimize resource utilization. However, the dynamic nature of network topology and resources poses challenges in establishing deterministic routing. These challenges encompass the intricacy of jointly scheduling transmission links and cycles, as well as the difficulty of maintaining stable end-to-end (E2E) routing paths. To tackle these challenges, our work introduces an efficient temporal graph-based deterministic routing strategy. Initially, we utilize a time-expanded graph (TEG) to represent the heterogeneous resources of an NTN in a time-slotted manner. With TEG, we meticulously define each necessary constraint and formulate the deterministic routing problem. Subsequently, we transform this nonlinear problem equivalently into solvable integer linear programming (ILP), providing a robust yet time-consuming performance upper bound. To address the considered problem with reduced complexity, we extend TEG by introducing virtual nodes and edges. This extension facilitates a uniform representation of heterogeneous network resources and traffic transmission requirements. Consequently, we propose a polynomial-time complexity algorithm, enabling the dynamic selection of optimal transmission links and cycles on a hop-by-hop basis. Simulation results validate that the proposed algorithm yields significant performance gains in traffic acceptance, justifying its additional complexity compared to existing routing strategies.

Index Terms—Non-terrestrial network, deterministic routing, temporal graph, integer linear programming, resource utilization.

I. INTRODUCTION

Non-terrestrial networks (NTNs) have emerged as a promising solution for global high-speed Internet access, thanks to their extensive coverage and robust bandwidth capabilities [1]. The continuous progress in aerial and space technologies, coupled with reduced manufacturing and launch costs, has notably hastened the development of NTNs. This acceleration is exemplified by the rapid establishment of mega-constellations like Starlink, OneWeb, and Telesat, highlighting the growing significance of NTNs in the future connectivity landscape. The 3rd Generation Partnership Project (3GPP) is actively dedicated to evolving 5G systems to support NTNs. Since 2017, 3GPP has released a series of documents focusing on network architecture, system configuration, and radio access [2]. Concurrently, the Internet Engineering Task Force (IETF) network working group diligently analyzes the requirements

of satellite constellations in the future Internet. Their analysis identifies efficient routing as a key enabler to enhance service performance and resource utilization within NTNs [3].

Nevertheless, designing efficient routing strategies for NTNs is challenging due to the dynamics of their network topologies and resources [4]. Over the years, researchers have explored and implemented diverse routing strategies tailored to specific NTNs. One widely adopted approach is the shortest path routing algorithm (SPR), designed to facilitate routing for pre-defined remote sensing transmission missions [5]. SPR models the network topology as a static graph over the mission duration, identifying the end-to-end (E2E) path with the minimum delay or hops. However, this approach lacks adaptability to changing network conditions and traffic demands. The snapshot graph-based routing algorithm (STR) extends SPR by employing a series of static snapshots to model the timevarying network topology and calculates E2E routing in each snapshot [6]. However, STR may not be able to determine feasible routing paths within a single snapshot when contacts or resources are scarce. Another routing strategy, the contact graph routing algorithm (CGR), incorporates the caching-andforwarding capability of satellites in routing decisions, enabling multi-hop transmissions in disruption-tolerant scenarios [7]. However, CGR prioritizes establishing the earliest connected E2E routing, potentially compromising optimal delay performance. Furthermore, the aforementioned strategies base routing decisions on bandwidth requirements over a long time duration, without allowing for the precise specification of traffic transmission times. Consequently, micro-bursts occur frequently, leading to uncertain delays and congestion.

Deterministic routing holds considerable promise within NTNs, offering the potential to enhance service performance and optimize resource utilization [8]. This technology facilitates precise scheduling of transmission links and cycles at each hop along the routing path, ensuring strict adherence to E2E delay and jitter requirements. Moreover, it enables dynamic allocation of network resources within each cycle on demand, thereby improving the overall resource utilization. Despite the commendable efforts of the Institute of Electrical and Electronics Engineers (IEEE) Time-Sensitive Networking (TSN) and the Internet Engineering Task Force (IETF) Deterministic Networking (DetNet) committees [9], the implementation of deterministic routing in NTNs remains challenging. This challenge stems from two primary factors: i) the high complexity associated with solving integer linear programming (ILP) for joint routing and scheduling falls short of meeting real-time processing requirements, and ii) the dynamic nature of network topology and resources complicates the identification of stable E2E routing paths.

K. Shi and H. Li are with the School of Communications Engineering, Xidian University, Xi'an 710071, P. R. China. (Email: kyshi10091@163.com; hyli@xidian.edu.cn).

J. Wang is with the Peng Cheng Laboratory, Shenzhen 518055, China (e-mail: wangjc61s@163.com).

K. Wang is with the School of Computer Science and Engineering, Xi'an University of Technology, Xi'an 710048, China. (Email: wangkan@xaut.edu.cn).

In response to these challenges, we introduce a temporal graph-based strategy to efficiently address the problem.

- Initially, we utilize a time-expanded graph (TEG) [10] to represent the heterogeneous resources of an NTN in a time-slotted manner, including contact topology, link capacity, node storage, link delay, and storage delay.
- With TEG, we formulate the deterministic routing problem comprehensively, incorporating a set of crucial constraints. Subsequently, we transform this nonlinear problem equivalently into a solvable ILP format. This transformation involves the linearization of cross-cycle propagation and caching constraints arising from long link delay and potential storage delay, thereby providing a robust yet time-consuming performance upper bound.
- To address the problem with reduced complexity, we construct an extended TEG (ETEG) to uniformly represent heterogeneous network resources and traffic transmission requirements. With ETEG, we propose a polynomial-time complexity algorithm for determining optimal transmission links and cycles on a hop-by-hop basis.
- Furthermore, we analyze the optimality and complexity of the proposed algorithm, followed by an implementation framework based on segment routing to facilitate its feasibility within large-scale NTNs. Simulation results demonstrate the superior performance of our proposal over SPR, STR, and CGR in terms of traffic acceptance. Additionally, it exhibits a significantly lower running time than the ILP-based strategy (referred to as ILPS).

II. SYSTEM MODEL AND PROBLEM FORMULATION

A. System model

In Fig. 1(a), we analyze an NTN, comprising N satellites denoted by the set $V = \{u_1, u_2, ..., u_N\}$. These satellites are interconnected through time-varying yet predictable transmission links. To enable precise delay control, the time window of each satellite is finely divided into consecutive cycles $\{\tau_h|h\geqslant 1\}$ of equal duration $|\tau|$, where $\tau_h=$ $((h-1)\cdot|\tau|,h\cdot|\tau|]$. This division allows for treating the network topology as static within each cycle. Additionally, we consider a time-critical (TC) traffic demand, denoted as f, characterized by a period of T_f and a per-period size of A_f . Assume that f is injected into the NTN at time t_f and needs to be delivered from a source satellite, $s \in V$, to a destination satellite, $d \in V$, within an upper bound of E2E delay, B_f . Then, we select the planning horizon, D, spanning from the start-cycle, $\tau_{\tilde{h}}$, to the end-cycle, $\tau_{\hat{h}}$, where t_f and t_f+B_f fall within $\tau_{\tilde{h}}$ and $\tau_{\hat{h}}$, respectively 1 . As illustrated in Fig. 1(c), we utilize a TEG, denoted as $\mathcal{G} = \{\mathcal{V}, \mathcal{E}, \mathcal{C}, \mathcal{L}\}\$, to model heterogeneous network resources in a time-slotted manner. Specifically, \mathcal{G} includes:

• A set of nodes, denoted as $\mathcal{V} = \left\{ u^h | u \in V, \check{h} \leqslant h \leqslant \widehat{h} \right\}$, where each $u^h \in \mathcal{V}$ signifies a satellite u within cycle τ_h .

 1 To enable deterministic transmission in all traffic periods, we establish deterministic routing for f in the first period and evolve it through repetition with cycle offset or necessary revisions.

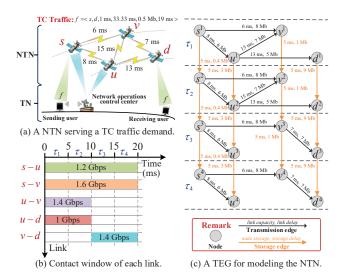


Fig. 1. Modeling a typical NTN using a TEG.

- A set of edges, denoted as \mathcal{E} , encompasses both transmission edges \mathcal{E}_t and storage edges \mathcal{E}_s . Herein, $\mathcal{E}_t = \left\{ (u^h, v^h) | u^h, v^h \in \mathcal{V}, \check{h} \leqslant h \leqslant \widehat{h} \right\}$ denotes the transmission links between satellites u and v in τ_h . Additionally, $\mathcal{E}_s = \left\{ (u^h, u^{h+1}) | u^h, u^{h+1} \in \mathcal{V}, \check{h} \leqslant h \leqslant \widehat{h} 1 \right\}$ depicts the capability of each satellite u to cache data across adjacent cycles (e.g., from τ_h to τ_{h+1}).
- A capacity set, denoted as \mathcal{C} , comprises two distinct subsets: a link capacity subset, \mathcal{C}_t , and a node storage subset, \mathcal{C}_s . Herein, $\mathcal{C}_t = \left\{ c_{u^h,v^h} | (u^h,v^h) \in \mathcal{E}_t, \check{h} \leqslant h \leqslant \widehat{h} \right\}$ signifies the maximum amount of data that can be transmitted on each transmission edge (u^h,v^h) , measured in megabytes (Mb). Furthermore, $\mathcal{C}_s = \left\{ c_{u^h,u^{h+1}} | (u^h,u^{h+1}) \in \mathcal{E}_s, \check{h} \leqslant h \leqslant \widehat{h} 1 \right\}$ represents the on-board storage resources of each satellite u during any cycle τ_h , measured in Mb.
- A delay set, denoted as \mathcal{L} , encompasses two distinct subsets: a link delay subset, \mathcal{L}_t , and and a storage delay subset, \mathcal{L}_s . Specifically, $\mathcal{L}_t = \left\{l_{u^h,v^h}|(u^h,v^h)\in\mathcal{E}_t,\check{h}\leqslant h\leqslant \hat{h}\right\}$ represents the propagation delay on each transmission edge (u^h,v^h) , measured in milliseconds (ms). Additionally, $\mathcal{L}_s = \left\{l_{u^h,u^{h+1}}|(u^h,u^{h+1})\in\mathcal{E}_s,\check{h}\leqslant h\leqslant \hat{h}-1\right\}$ depicts the cross-cycle caching delay (e.g., from τ_h to τ_{h+1}) at each satellite u, measured in ms. Without loss of generality, we can set $l_{u^h,u^{h+1}}=|\tau|$ for any $l_{u^h,u^{h+1}}\in\mathcal{L}_s$.

B. Constraint establishment

The deterministic routing problem focuses on effectively scheduling E2E transmission for the TC traffic demand f. This entails a judicious selection of satellites and links, along with the identification of suitable cycles for the transmission and caching of f at each satellite. To facilitate a comprehensive problem formulation, we define binary-valued variables for all edges in \mathcal{G} , denoted as $X = \left\{x_{u,v}^{h,k}|(u^h,v^k) \in \mathcal{E}\right\}$, accounting for two distinct cases: for a transmission edge $(u^h,v^h) \in \mathcal{E}_t$,

 $x_{u,v}^{h,h}=1$ indicates that f is transmitted from satellite u to satellite v within the cycle τ_h ; otherwise, $x_{u,v}^{h,h}=0$. Concerning a storage edge $(u^h,u^{h+1})\in\mathcal{E}_s$, if f is cached by satellite u from cycle τ_h to cycle τ_{h+1} , then $x_{u,u}^{h,h+1}=1$; otherwise, $x_{u,u}^{h,h+1}=0$.

1) E2E transmission constraint: f is required to initiate from the source satellite, s, and be transmitted to the destination satellite, d, within the planning horizon, D, i.e.,

$$\sum_{v^h: (s^h, v^h) \in \mathcal{E}_t} \sum_{h = \check{h}}^{\hat{h}} x_{s,v}^{h,h} = \sum_{u^h: (u^h, d^h) \in \mathcal{E}_t} \sum_{h = \check{h}}^{\hat{h}} x_{u,d}^{h,h} = 1.$$
 (1)

2) Lossless forwarding constraint: For any given satellite v (excluding s and d), it must forward the received f within D, expressed as:

$$\sum_{u^h:(u^h,v^h)\in\mathcal{E}_t}\sum_{h=\check{h}}^{\hat{h}}x_{u,v}^{h,h} = \sum_{w^h:(v^h,w^h)\in\mathcal{E}_t}\sum_{h=\check{h}}^{\hat{h}}x_{v,w}^{h,h}, \forall v\in V\setminus\{s,d\}. \quad (2)$$

3) Capacity constraint: For any transmission edge selected to transmit f, its link capacity should be no less than the traffic size A_f , expressed as:

$$c_{u^h,v^h} - x_{u,v}^{h,h} \cdot A_f \geqslant 0, \forall \left(u^h, v^h\right) \in \mathcal{E}_t. \tag{3}$$

4) Storage constraint: For any storage edge to cache f, its node storage should be no less than A_f , formulated as:

$$c_{u^h,u^{h+1}} - x_{u,u}^{h,h+1} \cdot A_f \geqslant 0, \forall \left(u^h, u^{h+1}\right) \in \mathcal{E}_s. \tag{4}$$

5) Cross-cycle propagation and caching constraints: When f is transmitted from satellite u to satellite v, the arrival time of f at v should occur later than the transmission cycle of u but no later than the transmission cycle of v, i.e.,

$$O(u)|\tau| - M \cdot (1 - y_{u,v}) < y_{u,v} \cdot [L_t(u) + L_s(u) + t_f] + l_{u,v} \le O(v)|\tau|$$
. (5a) Here,

$$O(u) = \sum_{(u^h, w^h) \in \mathcal{E}_t} \sum_{h=\tilde{h}}^{\hat{h}} h \cdot x_{u,w}^{h,h}$$
 (5b)

represents the transmission cycle² of u,

$$y_{u,v} = \sum_{h=\check{h}}^{\check{h}} x_{u,v}^{h,h},$$
 (5c)

indicates whether f is transmitted from u to v,

$$L_{t}(u) = \sum_{(p^{k}, q^{k}) \in \mathcal{E}_{t}} \sum_{k = \check{h}}^{O(u) - 1} x_{p, q}^{k, k} \cdot l_{p^{k}, q^{k}}$$
 (5d)

represents the total propagation delay from s to u,

$$L_s(u) = \sum_{(p^k, p^{k+1}) \in \mathcal{E}_s} \sum_{k=\tilde{h}}^{O(u)-1} x_{p,p}^{k,k+1} \cdot l_{p^k, p^{k+1}}$$
 (5e)

 2 If O(u) = 0, it indicates that u does not transmit f.

represents the total caching delay s to u, together with the caching delay within u,

$$l_{u,v} = \sum_{k=\check{b}}^{\hat{h}} x_{u,v}^{k,k} \cdot l_{u^k,v^k},$$
 (5f)

represents the propagation delay from u and v within the transmission cycle of u, and M represents a very large positive constant.

6) Transmission timing constraint: The time when satellite u sends f satellite u should fall within the transmission cycle of u, i.e.,

$$[O(u)-1]|\tau|-M\cdot(1-y_{u,v}) \leq y_{u,v}\cdot[L_t(u)+L_s(u)+t_f] \leq O(u)|\tau|,$$
 (6)

7) Caching timing constraint: When satellite v needs to cache f received from satellite u, the feasible cycles for caching should be no earlier that the cycle within which f arrives at v but earlier than the transmission cycle of v, expressed as:

$$y_{u,v} \cdot [L_t(u) + L_s(u) + t_f] + l_{u,v} - M \cdot (1 - x_{v,v}^{h,h+1}) \leqslant h \cdot x_{v,v}^{h,h+1} \cdot |\tau|$$

$$< O(v) \cdot |\tau| + \varepsilon \cdot (1 - x_{v,v}^{h,h+1}), \forall \check{h} \leqslant h \leqslant \hat{h}, \quad (7)$$

where ε represents a very small positive constant.

C. Problem formulation

$$\mathbf{P_{1}} : \min_{X} \sum_{(u^{h}, v^{h}) \in \mathcal{E}_{t}} \sum_{h = \check{h}}^{\check{h}} x_{u, v}^{h, h} \cdot l_{u^{h}, v^{h}} + \sum_{(u^{h}, u^{h+1}) \in \mathcal{E}_{s}} \sum_{h = \check{h}}^{\check{h} - 1} x_{u, u}^{h, h+1} \cdot l_{u^{h}, u^{h+1}}$$
(8)
s.t.(1) - (4), (5a), (6), (7).

The objective of $\mathbf{P_1}$ is to minimize the E2E delay, encompassing both the total propagation delay and caching delay during the delivery of f from s to d. Assuming the objective value is less than B_f , the solution to $\mathbf{P_1}$ provides efficient transmission scheduling for f with deterministic guarantees. However, due to the presence of the term $y_{u,v} \cdot [L_t(u) + L_s(u) + t_f]$ in constraints (5a), (6), and (7), where $L_t(u)$ and $L_s(u)$ involve variable-dependent summation and $y_{u,v}$ undergoes multivariable multiplication with both, these constraints become nonlinear. Consequently, $\mathbf{P_1}$ is unsolvable using existing ILP solvers. To address this challenge, we linearize these constraints by introducing auxiliary binary-valued variables along with a set of linear constraints.

Initially, we transform the variable-dependent summation in $L_t(u)$ and $L_s(u)$ into a summation term for the product of independent variables. This is achieved by introducing an auxiliary binary-valued variable, χ^u , and the following linear constraints:

$$k - O(u) \ge (2 - H) \cdot \chi^u + \varepsilon \cdot (1 - \chi^u) - 1, \forall \check{h} \le k \le \hat{h},$$
 (9)

and

$$k - O(u) \leq (H - 1) \cdot (1 - \chi^u) - \varepsilon \cdot \chi^u, \forall \check{h} \leq k \leq \hat{h}.$$
 (10)

Here, $H = \hat{h} - \check{h} + 1$. Using χ^u , $L_t(u)$ and $L_s(u)$ turn to be

$$L_t^{(1)}(u) = \sum_{(p^k, q^k) \in \mathcal{E}_t} \sum_{k=\tilde{h}}^h \chi^u \cdot x_{p,q}^{k,k} \cdot l_{p^k, q^k}, \tag{11}$$

and

$$L_s^{(1)}(u) = \sum_{(p^k, p^{k+1}) \in \mathcal{E}_s} \sum_{k=\tilde{h}}^{\hat{h}-1} \chi^u \cdot x_{p,p}^{k,k+1} \cdot l_{p^k, p^{k+1}}.$$
 (12)

Subsequently, we deal with multivariable multiplication in $y_{u,v} \cdot L_t^{(1)}(u)$ and $y_{u,v} \cdot L_t^{(1)}(u)$. For simplicity, we establish a general transformation paradigm for terms with the form $\prod_{m=1}^M a_m$, where $a_m \in \{0,1\}$. Specifically, we introduce an auxiliary binary-valued

variable, denoted as $\tilde{a} = \prod_{m=1}^{M} a_m$, for substitution, adhering to the following linear constraints:

$$\sum_{m=1}^{M} a_m - M + 1 \leqslant \widetilde{a} \leqslant a_m, \forall 1 \leqslant m \leqslant M.$$
 (13)

Based on the paradigm, we can respectively transform $y_{u,v}$. $L_t^{(1)}(u)$ and $y_{u,v} \cdot L_t^{(1)}(u)$ into

$$L_{t}^{(2)}(u) = \sum_{(p^{k}, q^{k}) \in \mathcal{E}_{t}} \sum_{h, k = \check{h}}^{\hat{h}} \widetilde{x}_{u, v, p, q}^{h, k} \cdot l_{p^{k}, q^{k}}$$
(14)

and

$$L_s^{(2)}(u) = \sum_{(p^k, p^{k+1}) \in \mathcal{E}_s} \sum_{k=\check{h}h=\check{h}}^{\hat{h}-1} \sum_{u,v,p,p}^{\hat{h}} \widetilde{x}_{u,v,p,p}^{h,k} \cdot l_{p^k,p^{k+1}},$$
(15)

with the introduced variables $\widetilde{x}_{u,v,p,p}^{h,k}$ and $\widetilde{x}_{u,v,p,p}^{h,k}$ satisfying

$$x_{u,v}^{h,h} + x_{p,q}^{k,k} + \chi^{u} - 2 \leqslant \widetilde{x}_{u,v,p,q}^{h,k} \leqslant x_{u,v}^{h,h}, \tag{16}$$

$$x_{u,v}^{h,h} + x_{p,q}^{k,k} + \chi^{u} - 2 \leqslant \widetilde{x}_{u,v,p,q}^{h,k} \leqslant x_{p,q}^{k,k}, \tag{17}$$

$$x_{u,n}^{h,h} + x_{n,a}^{k,k} + \chi^u - 2 \leqslant \tilde{x}_{u,n,a}^{h,k} \leqslant \chi^u,$$
 (18)

$$x_{u,v}^{h,h} + x_{p,p}^{k,k+1} + \chi^{u} - 2 \leqslant \widetilde{x}_{u,v,p,p}^{h,k} \leqslant x_{u,v}^{h,h}, \tag{19}$$

$$x_{u,v} + x_{p,q} + \chi - 2 \leqslant x_{u,v,p,q} \leqslant x_{p,q},$$

$$x_{u,v}^{h,h} + x_{p,q}^{k,k} + \chi^{u} - 2 \leqslant \widetilde{x}_{u,v,p,q}^{h,k} \leqslant \chi^{u},$$

$$x_{u,v}^{h,h} + x_{p,p}^{k,k+1} + \chi^{u} - 2 \leqslant \widetilde{x}_{u,v,p,p}^{h,k} \leqslant x_{u,v}^{h,h},$$

$$x_{u,v}^{h,h} + x_{p,p}^{k,k+1} + \chi^{u} - 2 \leqslant \widetilde{x}_{u,v,p,p}^{h,k} \leqslant x_{p,p}^{h,h},$$

$$x_{u,v}^{h,h} + x_{p,p}^{k,k+1} + \chi^{u} - 2 \leqslant \widetilde{x}_{u,v,p,p}^{h,k} \leqslant x_{p,p}^{k,k+1},$$

$$x_{u,v}^{h,h} + x_{p,p}^{k,k+1} + \chi^{u} - 2 \leqslant \widetilde{x}_{u,v,p,p}^{h,k} \leqslant \chi^{u}.$$

$$(21)$$

$$x_{u,v}^{h,h} + x_{n,n}^{k,k+1} + \chi^u - 2 \leqslant \widetilde{x}_{u,v,n,n}^{h,k} \leqslant \chi^u. \tag{21}$$

Substituting (14) and (15) into (5a), (6), and (7), we can obtain the following linear constraints:

$$O(u)|\tau| - M \cdot (1 - y_{u,v}) < L_t^{(2)}(u) + L_s^{(2)}(u) + y_{u,v} \cdot t_f + l_{u,v} \le O(v)|\tau|,$$
 (22)

$$[O(u)-1] |\tau| - M \cdot (1-y_{u,v}) \leqslant L_t^{(2)}(u) + L_s^{(2)}(u) + y_{u,v} \cdot t_f \leqslant O(u) |\tau|, \ \ (23)$$

$$\begin{split} L_t^{(2)}(u) + L_s^{(2)}(u) + y_{u,v} \cdot t_f + l_{u,v} - M \cdot (1 - x_{v,v}^{h,h+1}) &\leqslant h \cdot x_{v,v}^{h,h+1} \cdot |\tau| \\ &< O(v) \cdot |\tau| + \varepsilon \cdot (1 - x_{v,v}^{h,h+1}), \, \forall \check{h} \leqslant h \leqslant \hat{h}. \end{split} \tag{24}$$

Ultimately, P_1 can be reformulated as an ILP problem:

$$\mathbf{P_2} : \min_{X'} \sum_{(u^h, v^h) \in \mathcal{E}_t} \sum_{h = \check{h}}^{\hat{h}} x_{u, v}^{h, h} \cdot l_{u^h, v^h} + \sum_{(u^h, u^{h+1}) \in \mathcal{E}_s} \sum_{h = \check{h}}^{\hat{h} - 1} x_{u, u}^{h, h+1} \cdot l_{u^h, u^{h+1}} \quad (8)$$
s.t. $(1) - (4), (9), (10), (16) - (24),$

where the set of decision variables is defined as $X' = X \bigcup \{\chi^u | u \in V\} \bigcup \{\widetilde{x}_{u,v,p,q}^{h,k}| (u^h,v^h), (p^k,q^k) \in \mathcal{E}_t\} \bigcup \{\widetilde{x}_{u,v,p,p}^{h,k}| (u^h,v^h) \in \mathcal{E}_t, (p^k,p^{k+1}) \in \mathcal{E}_s\}$. Notably, ILP solvers [11] can effective formula of the set of decision variables is defined as $X' = X \cup \{\chi^u | u \in V\} \cup \{\widetilde{x}_{u,v,p,p}^{h,k} | u \in V\} \cup \{\widetilde{x}_{$ tively handle P2, providing a robust performance upper bound for P_1 . Nevertheless, the computations involved in these solvers prove excessively time-consuming, failing to meet real-time processing requirements. Therefore, we introduce a graph-based method to address P_1 , aiming to reduce the running time while ensuring optimality.

III. TEMPORAL GRAPH-BASED DETERMINISTIC ROUTING

A. Extended time-expanded graph model

To effectively address P_1 , we enhance the original graph \mathcal{G} to form an ETEG, denoted as $\mathcal{G}' = \{\mathcal{V}', \mathcal{E}', \mathcal{C}', \mathcal{L}'\}$, by introducing virtual nodes and edges to establish a uniform representation of heterogeneous network resources and traffic transmission requirements. The construction of \mathcal{G}' is outlined as follows, as illustrated in Fig. 2:

- 1) To capture time-slotted network resources, we initialize the node set \mathcal{V}' as \mathcal{V} , the edge set \mathcal{E}' as \mathcal{E} , the capacity set \mathcal{C}' as \mathcal{C} , and the delay set \mathcal{L}' as \mathcal{L} .
- 2) To signify the earliest cycle within which f departs from the source satellite s, we introduce a virtual source s' into \mathcal{V}' and a virtual transmission edge (s', s^h) into \mathcal{E}' .

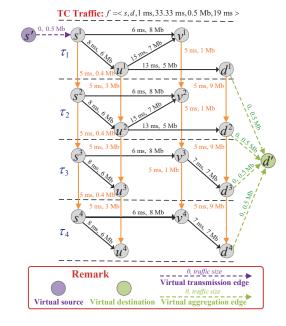


Fig. 2. An ETEG model.

- 3) To represent potential cycles within which f is transmitted to the destination satellite d, we introduce a virtual destination d' into \mathcal{V}' and a set of virtual aggregation edges $\{(d^h, d') | \check{h} \leq h \leq \widehat{h}\}$ into \mathcal{E}' . By designating d' as the unique destination, we can avoid traversing potential routing to the destination satellite within each cycle (i.e., d^h , where $h \leq h \leq h$), thus significantly reducing the computational complexity of deterministic routing.
- 4) To indicate the capacity requirement of f, we set the capacity metrics of all introduced virtual transmission edges and virtual aggregation edges as A_f . Since these edges lack physical counterparts, we assign a delay metric of 0 to them, thus not affecting the overall E2E delay.

B. ETEG-based deterministic routing algorithm

Based on ETEG, we identify that the deterministic routing problem is equivalent to a path-finding one, rather than directly solving the ILP in P2. Using this insight, we propose an ETEG-based deterministic routing algorithm with low complexity (as detailed in Algorithm 1). The proposed algorithm jointly utilizes both link capacity and node storage, facilitating cross-cycle propagation and caching of f. Consequently, it dynamically selects optimal links and cycles on a hop-by-hop basis to establish a time-featured path (Definition 1) that minimizes the E2E delay while meeting resource requirements.

Definition 1. A time-featured path, denoted as \mathcal{P}_f , can be represented by a node sequence $s' \to \cdots \to u^h \to v^{k'} \to \cdots \to d'$ in the ETEG, adhering to the condition: if $u \neq v$, then h < kand $c_{u^h,v^h} \ge A_f$; otherwise k = h + 1 and $c_{u^h,v^k} \ge A_f$.

In Algorithm 1, we execute the path-finding process based on $\mathcal{G}' = \{\mathcal{V}', \mathcal{E}', \mathcal{C}', \mathcal{L}'\}$. To facilitate this process, we define two essential parameters at each node $u^h \in \mathcal{V}'$: the prenode $p(u^h)$, indicating the previous hop of u^h in the timefeatured path \mathcal{P}_f , and the node delay $L(u^h)$, representing the propagation delay and caching delay along \mathcal{P}_f from s'to u^h . Additionally, we introduce a priority queue, denoted as Q, for maintaining nodes awaiting the determination of

Algorithm 1 ETEG-based deterministic routing algorithm

```
1: Input: \mathcal{G}' = \{\mathcal{V}', \mathcal{E}', \mathcal{C}', \mathcal{L}'\}.
 2: \mathcal{P}_f = \varnothing, p(u^h) = \varnothing, \forall u^h \in \mathcal{V}', L(s') = t_f, L(v^k) = t_f
        +\infty, \forall v^k \in \mathcal{V}' \setminus \{s'\}, Q = \mathcal{V}';
 3: while d' \in Q do
               u^h = \arg\min L\left(v^k\right);
 4:
               Q \leftarrow Q \backslash \{u^h\};
 5:
               for v^k \in Q: (u^h, v^k) \in \mathcal{E}' do
 6:
                       if c_{u^h,v^k} \geqslant A_f then
 7:
                               \begin{split} r &= \left\lceil \frac{1}{|\tau|} \cdot \left( L\left(u^h\right) + l_{u^h,v^k}\right) \right\rceil; \\ \textbf{if } L\left(u^h\right) + l_{u^h,v^k} < L\left(v^r\right) \textbf{ then} \\ L\left(v^r\right) &= L\left(u^h\right) + l_{u^h,v^k}; \end{split} 
 8:
 9:
10:
11:
                               end if
12:
                       end if
13:
               end for
14:
15: end while
16: if L(d') \leq t_f + B_f then
               u^h = d':
17:
               while u^h \neq \emptyset do
18:
                       \mathcal{P}_f \leftarrow \mathcal{P}_f \bigcup \{u^h\};u^h = p(u^h);
19:
20:
21:
22: end if
23: Output: \mathcal{P}_f with deterministic guarantees.
```

their node delay. During initialization (in step 2), we set $\mathcal{P}_f = \varnothing$, $p\left(u^h\right) = \varnothing$ for any node $u^h \in \mathcal{V}'$, $L\left(s'\right) = t_f$, $L\left(v^k\right) = +\infty$ for any node $v^k \in \mathcal{V}'\backslash \{s'\}$, and $Q = \mathcal{V}'$. In each iteration (from steps 3 to 15), we extract u^h with the minimum node delay from Q. Subsequently, we update each node v^k adjacent to u^h , provided that the resources (i.e., link capacity or node storage) are sufficient, and the node delay of v^k can be reduced via the relay by u^h . Due to cross-cycle propagation and caching constraints, the updated node might be v^r with $r = \left\lceil \frac{1}{|\tau|} \cdot \left(L\left(u^h\right) + l_{u^h,v^k}\right) \right\rceil$. Consequently, we have $L\left(v^r\right) = L\left(u^h\right) + l_{u^h,v^k}$ and designate $p\left(v^r\right) = u^h$. The above iteration continues until d' is extracted from Q. Finally, if $L\left(d'\right) \leqslant t_f + B_f$ holds, we can obtain \mathcal{P}_f by backtracking from d' to s' (from steps 17 to 21); otherwise, no feasible \mathcal{P}_f exists.

Fig. 3 illustrates an application. Initiated from s', the proposed algorithm updates the node delay of its sole neighbor, s^1 , to $L\left(s^1\right)=1$ ms and set $p(s^1)=s'$. Next, we pop s^1 from $\mathcal Q$ and traverse all its neighbors, yielding the following: $L\left(s^2\right)=+\infty$, due to insufficient node storage of (s^1,s^2) ; $L\left(u^2\right)=9$ ms and $L\left(v^2\right)=7$ ms, owing to cross-cycle propagation from s^1 to u^2 and v^2 , respectively. In iteration 3, v^2 is extracted from Q, and $L\left(v^3\right)$ is updated to 7+5=12 ms through cross-cycle caching from v^2 to v^3 . The proposed algorithm then selects u^2 and update its sole neighbor u^3 with $L\left(u^3\right)=9+5=14$ ms. Subsequent steps involve updating $L\left(d^4\right)$ to 12+7=19 ms through cross-cycle propagation from v^3 to d^4 , which also serves as the ultimate node delay at d', i.e., $L\left(d'\right)=19$ ms <1+19=20 ms. Through

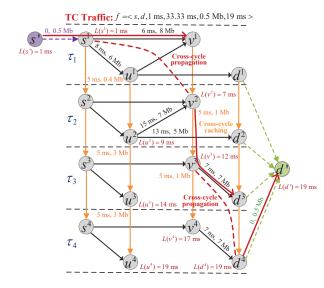


Fig. 3. An application of the ETEG-based deterministic routing algorithm.

backtracking, a feasible time-featured path $\mathcal{P}_f: s' \to s^1 \to v^2 \to v^3 \to d^4 \to d'$ is obtained.

C. Complexity and optimality analyses

Theorem 1. The ETEG-based deterministic routing algorithm is capable of calculating a time-featured path with minimum E2E delay.

Proof. For **Algorithm 1**, we demonstrate that each node extracted from Q has been determined with the minimum node delay. This assertion remains valid for s^1 with $L\left(s^1\right)=t_f$. Moving on to the M-th extracted node, denoted as u^h , with a node delay $L\left(u^h\right)$, we evaluate whether $L\left(u^h\right)$ can be further reduced through relaying by any node in Q, such as w^r . If so, either $L\left(w^r\right)+l_{w^r,u^r}< L\left(u^h\right)$ holds for cross-cycle propagation, or $L\left(w^r\right)+l_{w^r,w^{r+1}}< L\left(w^{r+1}\right)=L\left(u^h\right)$ holds for cross-cycle caching. However, due to step 4, $L\left(w^r\right)\geqslant L\left(u^h\right)$ holds, along with $l_{w^r,u^r}\geqslant 0$ and $l_{w^r,w^{r+1}}\geqslant 0$, thereby contradicting the aforementioned inequalities. Conversely, $L\left(u^h\right)$ is minimized when u^h is extracted from Q, also holding for d'. The proof is completed. \square

Theorem 2. The time complexity of the ETEG-based deterministic routing algorithm is $O(|\mathcal{E}'| \cdot \log |\mathcal{V}'|)$, where $|\mathcal{V}'|$ and $|\mathcal{E}'|$ denote the number of nodes and edges in the input \mathcal{G}' , respectively.

Proof. For **Algorithm 1**, assume that the \mathcal{G}' and Q are stored in adjacency lists and binary heaps, respectively. The initialization in step 2 takes $O(|\mathcal{V}'|)$ time. During each iteration from steps 3 to 15, it requires O(1) time to extract u^h with the minimum node delay from Q and $O(\log |\mathcal{V}'|)$ time to update Q. Furthermore, updating all nodes adjacent to u^h takes at most $O(\deg(u^h) \cdot \log |\mathcal{V}'|)$, where $\deg(u^h)$ represents the out-degree of u^h in \mathcal{G}' . At worst, we must traverse all nodes in \mathcal{V}' once before extracting d' from Q. Therefore, the time complexity reaches $O(|\mathcal{V}'| \cdot (1 + \log |\mathcal{V}'| + \deg(u^h) \cdot \log |\mathcal{V}'|)) =$

 $O\left((|\mathcal{V}'|+|\mathcal{E}'|)\cdot \log |\mathcal{V}'|\right)$. Additionally, the backtracking process takes at most $O\left(|\mathcal{V}'|\right)$. Thus, the total time complexity is $O\left(|\mathcal{V}'|\right)+O\left((|\mathcal{V}'|+|\mathcal{E}'|)\cdot \log |\mathcal{V}'|\right)+O\left(|\mathcal{V}'|\right)=O\left(|\mathcal{E}'|\cdot \log |\mathcal{V}'|\right)$, as $|\mathcal{E}'|\geqslant |\mathcal{V}'|-1$ for a connected graph \mathcal{G}' . The proof is completed.

D. Algorithm implementation

Following [12], we propose an implementation framework based on segment routing [13] for our algorithm. Within the NTN scenario in Fig. 1(a), we present the key aspects of the framework as follows:

- 1) Parameter maintenance: The network operations control center (NOCC) continuously acquires link status information from satellites through low propagation delay satellite-to-ground links. It extracts essential network parameters for deterministic routing decisions, including link capacity, link delay, and node storage.
- 2) Routing decision: Leveraging the maintained network parameters and traffic information from the terrestrial sending user, the NOCC constructs the ETEG and determines optimal deterministic routing for the TC traffic demand. Additionally, the NOCC not only reserves the required resources for the demand by updating network parameters but also configures the deterministic forwarding table sent to the specified source satellite.
- 3) Routing deployment: Following the deterministic forwarding table, the source satellite modifies TC traffic demand packets injected by the sending user. This modification involves encapsulating per-hop transmission link and cycle information into the packets' headers, directing them across the network transmission until reaching the destination satellite. Then, the packets are decapsulated and downlinked to the terrestrial receiving user.
- 4) Routing evolution: The NOCC consistently evaluates the feasibility of preceding deterministic routing for upcoming traffic periods. If feasible, the source satellite simply introduces a period-size cycle offset when deploying the deterministic forwarding table; otherwise, the NOCC re-executes the routing decision in 2) and the routing deployment in 3).

Notably, a cross-domain decision architecture [14] can be alternatively deployed when a single NOCC is insufficient to handle all TC traffic demands across the NTN, or when long propagation delays emerge as a primary concern.

IV. SIMULATIONS

A. Simulation setup

We conduct simulations on a partial Starlink constellation comprising 168 satellites selected from S1 [15]. These satellites are distributed across 12 orbits, each accommodating 14 satellites, positioned at a height of 550 km with an inclination of 53°. Employing the Satellite Toolkit (STK) simulator, we generate a time-varying NTN scenario with parameters in Table I. Consider that TC traffic demands continuously enter the NTN within the initial 120 seconds (s), following a

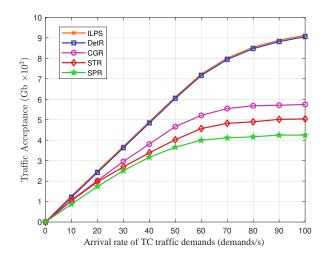


Fig. 4. Evaluation of traffic acceptance.

Poisson process. The source-destination satellite pairs of these demands are randomly specified across the constellation. Each demand, representing applications like high-quality video telephony [16], operates with a period of 33.33 ms (equivalent to a video frame rate of 30 frames per second) and has an active duration varying from 60 to 180 s. Furthermore, the per-cycle size of each demand follows a uniform distribution between 0.05 Mb and 0.6 Mb, with an E2E delay upper bound set at 75 ms. Simulations are executed on a Hewlett-Packard Z620 tower workstation (Intel Core i9-13900H CPU@2.60GHz, 32G RAM, Windows 11 x64) with a C++ environment.

TABLE I NETWORK PARAMETERS

Cycle duration (ms)	Link capacity (Mb)	Node storage (Mb)	Link delay (ms)
$ \tau = 5$	$c_{u^h,v^h} = 5$	$c_{u^h,u^{h+1}} = 1000$	$l_{v^h,v^h} \in [5,12]$

Throughout the entire 300-second horizon, we evaluate the performance of our proposed algorithm (referred to as DetR) and four benchmark strategies: SPR, STR, CGR, and ILPS, using the following metrics:

- Traffic acceptance α: The total size of demands with E2E deterministic transmission guarantees, providing insights into network resource utilization.
- Average running time β : The average time to process a single demand.
- Average E2E path delay γ : The average delay of routing paths associated with demands possessing E2E deterministic transmission guarantees.

B. Simulation results

We first evaluate the traffic acceptance, α , by varying the arrival rate from 1 to 100 demands per second (demands/s), as shown in Fig. 4. As expected, α values for all algorithms increase as the arrival rate increases, but the increase rate gradually slows down. This happens as the increasing demands occupy most of the available network resources. Notably, DetR

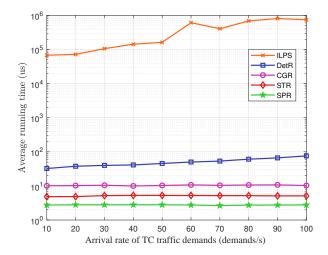


Fig. 5. Evaluation of average running time.

is comparable to the optimal ILPS and surpasses CGR, STR, and SPR. In particular, when the arrival rate reaches 100 demands/s, α is improved by more than 50%. This significant enhancement is attributed to the ability of DetR and ILPS to jointly utilize link capacity and node storage in different cycles, enabling conflict-free routing paths for a larger number of demands. In contrast, SPR records the lowest performance due to its neglect of time-varying network characteristics, focusing solely on routing paths within a static graph. Consequently, these paths may encounter interruptions and congestion. Notably, STR employs a series of time-evolving snapshots to model the NTN, while CGR introduces connectivity between adjacent snapshots, expanding the solution space compared to SPR. However, since the routing paths are determined based on average bandwidth requirements, micro-bursts occur frequently among demands, limiting traffic acceptance.

Fig. 5 illustrates the average running time, β , under varying arrival rates of demands. Notably, β values for four graphbased algorithms are significantly lower than that of ILPS, primarily attributed to not traversing the entire solution space for routing decisions. DetR is the slowest among the four, with the gap not exceeding 80 microseconds (us). The increased complexity arises because DetR determines optimal transmission links and cycles for demands within a time-expanded solution space. Together with Fig. 4, it becomes evident that the commendable enhancement achieved by DetR in traffic acceptance justifies its increased complexity compared to SPR, STR, and CGR. Furthermore, DetR's β exhibits a gradual increase with rising arrival rates since DetR facilitates deterministic routing evolution by introducing cycle offsets or performing re-execution. These actions become more frequent as demands increase, thus increasing β .

We also evaluate the average E2E path delay, γ , for both DetR and ILPS. As shown in Fig. 6, the γ of DetR aligns with that of ILPS, gradually increasing as the arrival rate of demands rises. This trend is expected, as network resources become limited under heavy demand loads. To accommodate more demands while maintaining conflict-free deterministic

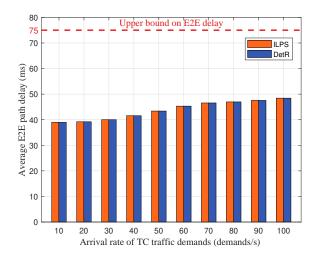


Fig. 6. Evaluation of average E2E path delay.

routing paths, essential cross-cycle propagation and caching are included in these paths. Although γ increases, it remains below the upper bound of 75 ms. Notably, the average E2E path delay of CGR, STR, and SPR is not presented, since their traffic acceptance is far lower than that of DetR, enabling them to find routing paths with lower delay under lightly loaded network conditions. Consequently, direct numerical comparisons among all algorithms lack meaningful insight.

V. CONCLUSION

This study focuses on addressing the deterministic routing problem within NTNs. Leveraging the TEG, we meticulously formulated and equivalently transformed this intricate problem into a solvable ILP problem, providing a robust yet time-consuming performance upper bound. To enhance time efficiency, we introduced an ETEG-based deterministic routing algorithm with polynomial time complexity. This algorithm enables the joint utilization of link capacity and node resources, facilitating cross-cycle propagation and caching of the TC traffic demands. Consequently, it can determine optimal transmission links and cycles on a hop-by-hop basis. Simulation results demonstrated that our proposal outperforms SPR, STR, and CGR in terms of traffic acceptance, thereby justifying its additional complexity. Furthermore, it exhibits significantly reduced running time compared to ILPS.

REFERENCES

- [1] M. Giordani et al., "Non-terrestrial Networks in the 6G Era: Challenges and Opportunities," *IEEE Netw.*, vol. 35, no. 2, pp. 244–251, Apr. 2021.
- [2] X. Lin et al., "5G from Space: An Overview of 3GPP Non-terrestrial Networks," *IEEE Commun. Stds. Mag.*, vol. 5, no. 4, pp. 147–153, 2021.
- [3] L. Han et al., "Problems and Requirements of Satellite Constellation for Internet," Internet Engineering Task Force, Internet-Draft draft-lhanproblems-requirements-satellite-net-03, Jul. 2022.
- [4] P. Wang et al., "Enhancing Earth Observation Throughput Using Intersatellite Communication," *IEEE Trans. Wireless Commun.*, vol. 21, no. 10, pp. 7990—8006, 2022.
- [5] T. Pan et al., "OPSPF: Orbit Prediction Shortest Path First Routing for Resilient LEO Satellite Networks," In *Proc. IEEE Int. Conf. Commun.*, 2019, pp. 1—6.

- [6] M. Werner, "A Dynamic Routing Concept for ATM-based Satellite Personal Communication Networks," *IEEE J. Sel. Areas Commun.*, vol. 15, no. 8, pp. 1636–1648, Oct. 1997.
- [7] G. Araniti et al., "Contact Graph Routing in NTN Space Networks: Overview, Enhancements and Performance," *IEEE Commun. Mag.*, vol. 53, no. 3, pp. 38–46, Mar. 2015.
 [8] B. Liu et al., "Towards Large-scale Deterministic IP Networks," In *IFIP*
- [8] B. Liu et al., "Towards Large-scale Deterministic IP Networks," In IFIP Networking Conf., Jun. 2021, pp. 1–9.
- [9] B. Varga et al., "Robustness and Reliability Provided by Deterministic Packet Networks (TSN and DetNet)," *IEEE Trans. Netw. Service Manag.*, early access, doi: 10.1109/TNSM.2023.3284590.
- [10] K. Shi et al., "Time-expanded Graph-based Energy-efficient Delay-bounded Multicast over Satellite Networks," *IEEE Trans. Veh. Technol.*, vol. 69, no. 9, pp. 10380–10384, Apr. 2020.
- [11] I. Gurobi Optimization, "Gurobi Optimizer Reference Manual," 2016. [Online]. Available: http://www.gurobi.com.
- [12] M. Hu, et al., "Software-defined Multicast Using Segment Routing in LEO Satellite Networks," *IEEE Trans. Mob. Comput.*, vol. 23, no. 1, pp. 835–849, Jan. 2024.
- [13] Z. N. Abdullah, et al., "Segment Routing in Software-defined Networks: A Survey," *IEEE Commun. Surv. Tutor.*, vol. 21, no. 1, pp. 464–486, Firstquarter 2019.
- [14] Y. Shi, et al., "A Cross-domain SDN Architecture for Multi-layered Space-terrestrial Integrated Networks," *IEEE Netw.*, vol. 33, no. 1, pp. 29–35, 2019.
- [15] A. U. Chaudhry, et al., "Laser Intersatellite Links in a Starlink Constellation: A Classification and Analysis," *IEEE Veh. Technol. Mag.*, vol. 16, no. 2, pp. 48–56, Jun. 2021.
- [16] H. Zhang et al., "OnRL: Improving Mobile Video Telephony via Online Reinforcement Learning," In *Proc. 26th Annu. Int. Conf. Mobile Comput. Netw.*, Sep. 2020, pp. 1–14.