

Faster SRv6 D-plane with XDP

Ryoga Saito

自己紹介

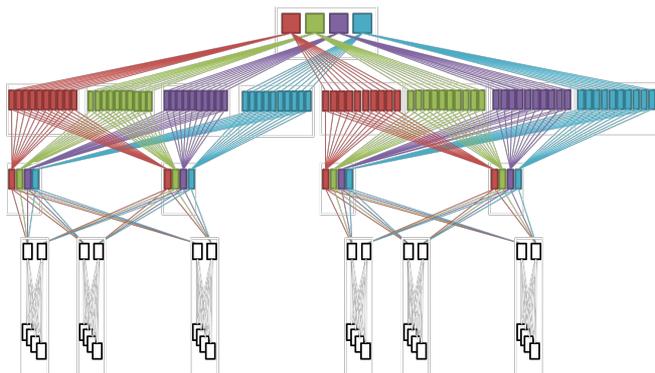
- 齋藤遼河
- 国立木更津工業高等専門学校 情報工学科 5年
- LINE株式会社
 - ITサービスセンター Verda室 ネットワーク開発チーム

Background

背景

Full L3 CLOS Network

- シングルテナントなネットワーク
- LINEのメッセージャーやファミリーサービスなどが稼働



特定サービス専用のネットワーク

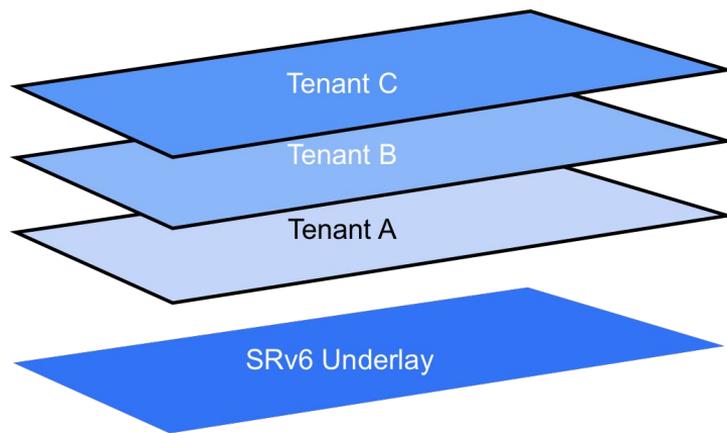
- サービス固有の要件を持つサービスが稼働
- サービス毎にネットワークを構築



アンダーレイネットワークの断片化
設計・構築にかかる時間の増加
運用コストの増加

背景

- アンダーレイネットワークを共通化し、運用負荷を軽減する
- オーバレイネットワークでテナント間分離し、それぞれのポリシーを実現する



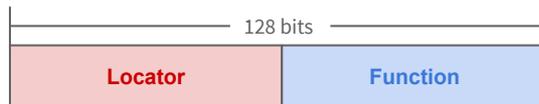
柔軟にスケールするオーバレイ
テナント毎に分離・独立したセキュリティ
将来的なサービスチェイニング

シンプルなL3のアンダーレイ

SRv6の概要

- Segment ID (SID)

- Segment (パケットに対する指示)を表す ID
 - **Locator**: Functionが実行されるノードを示す情報
 - **Function**: Functionを示す情報



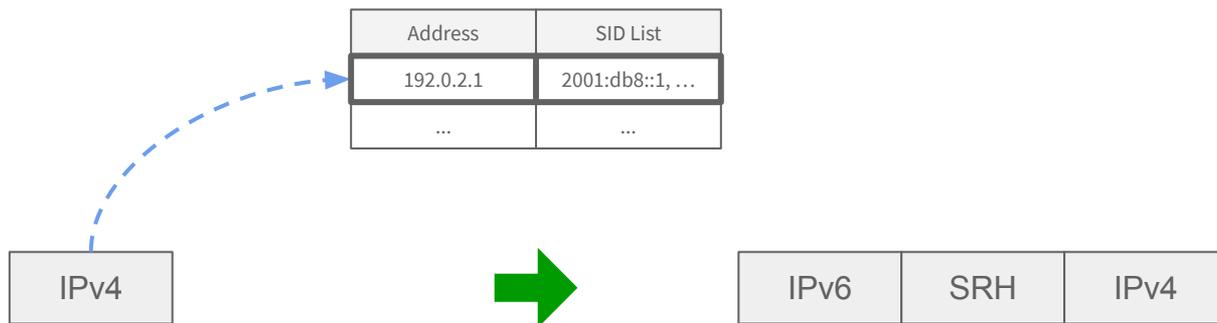
- Segment Routing Header (SRH)

- IPv6のルーティング拡張ヘッダ
- Segment List (通るべきSIDのリスト) や Segments Left (Segment List内の現在のSIDのオフセット) などが含まれている

SRv6の概要

- Well-known Functions

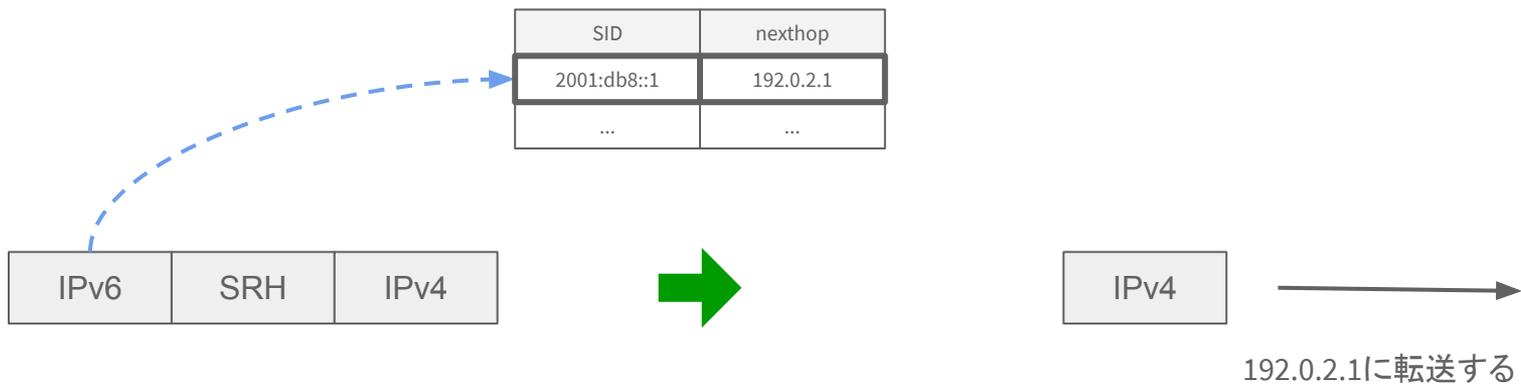
- H.Encaps (SR Headend with Encapsulation in an SRv6 Policy):
IPv6ヘッダ, SRHを用いてカプセル化し転送する



SRv6の概要

- Well-known Functions

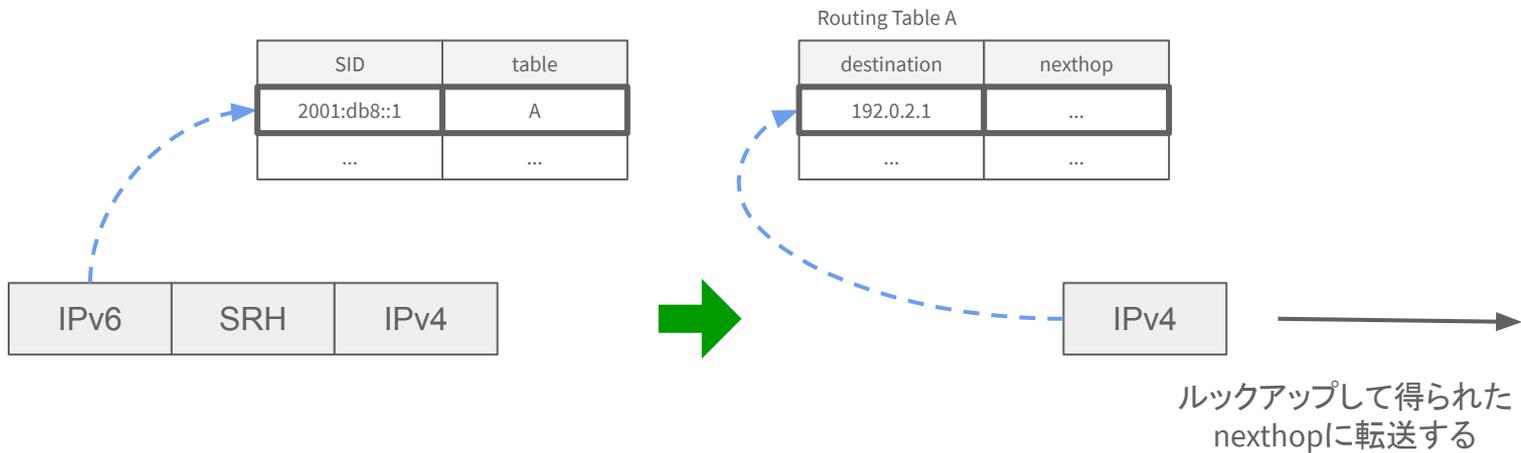
- End.DX4 (Decapsulation and IPv4 Cross-Connect):
IPv6ヘッダ, SRHを外し、指定された nexthopに転送する



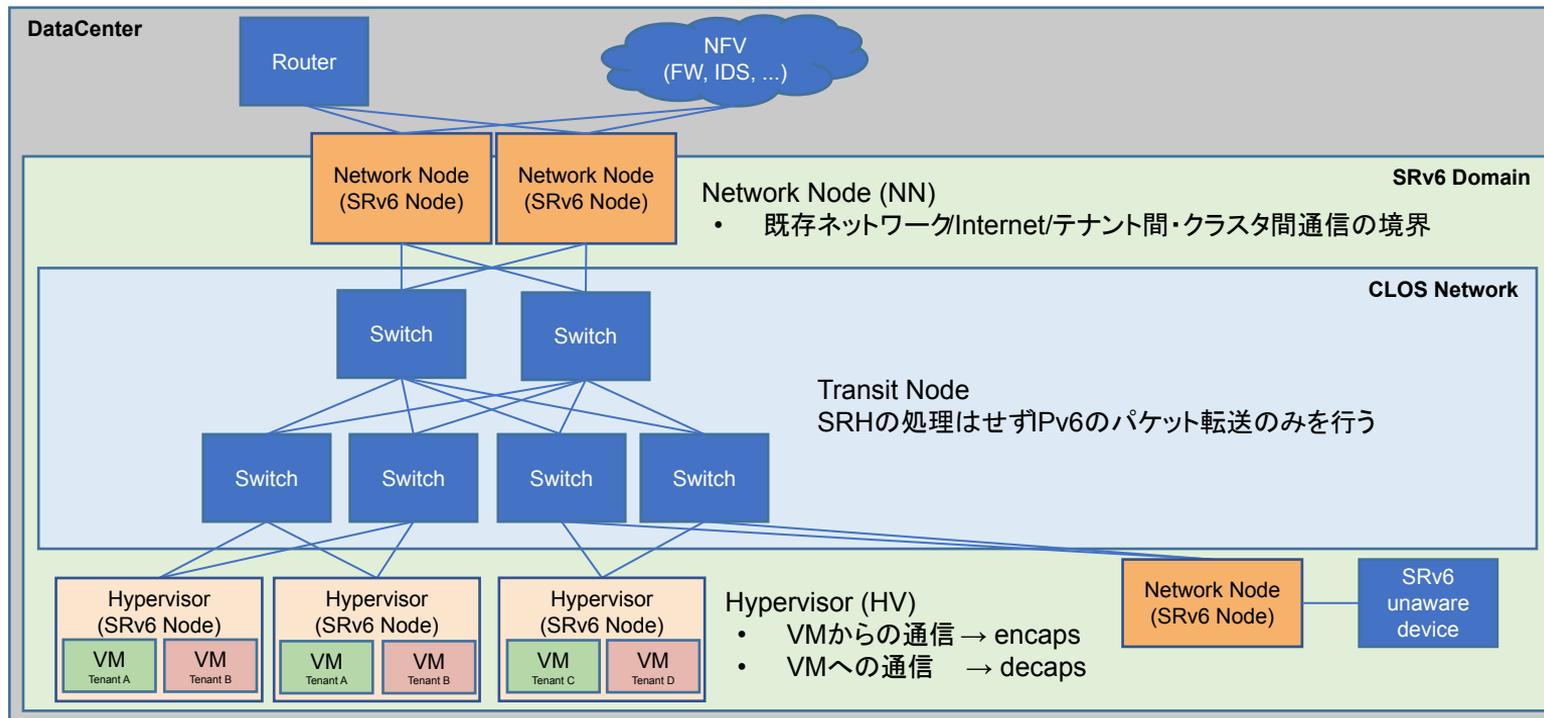
SRv6の概要

- Well-known Functions

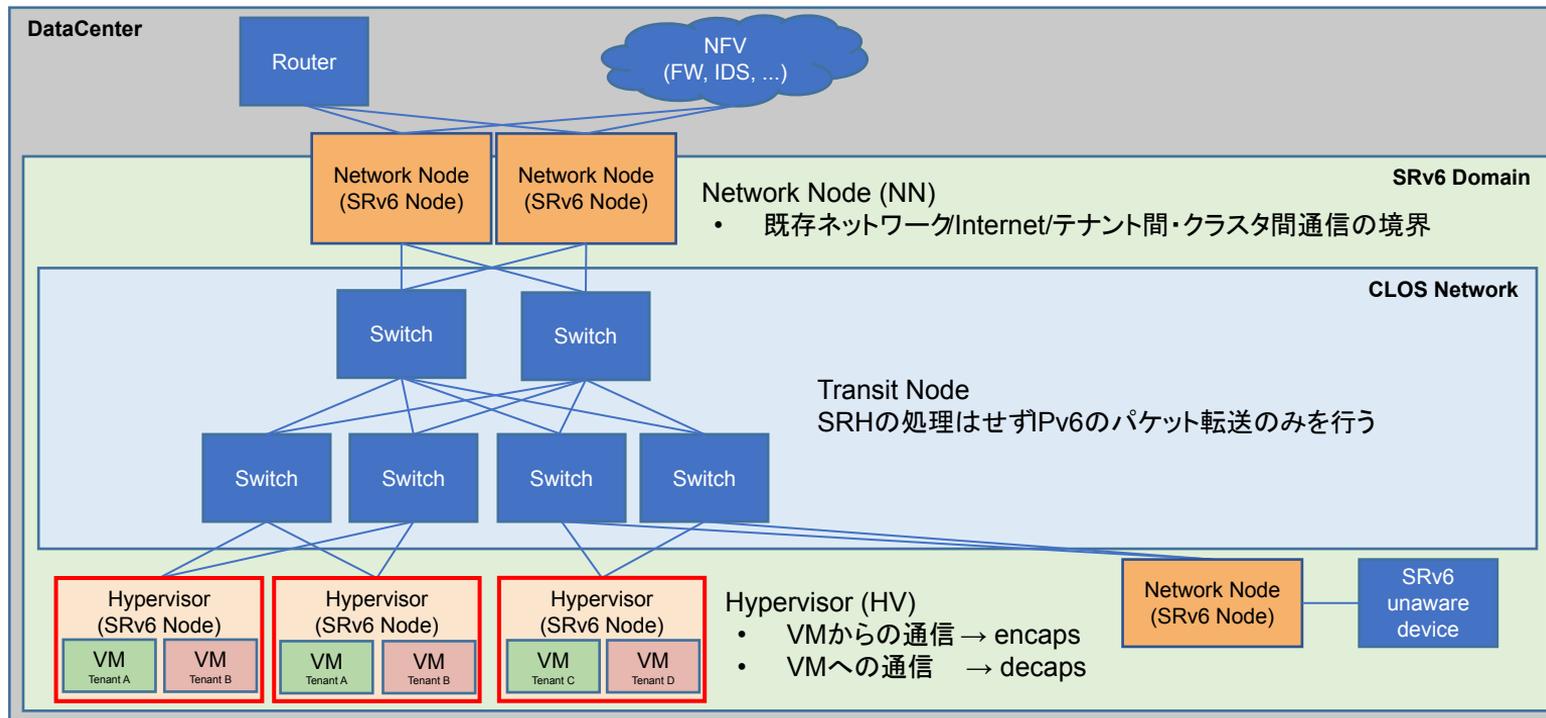
- End.DT4 (Decapsulation and Specific IPv4 Table Lookup):
IPv6ヘッダ, SRHを外し、指定されたルーティングテーブルを用いて転送する



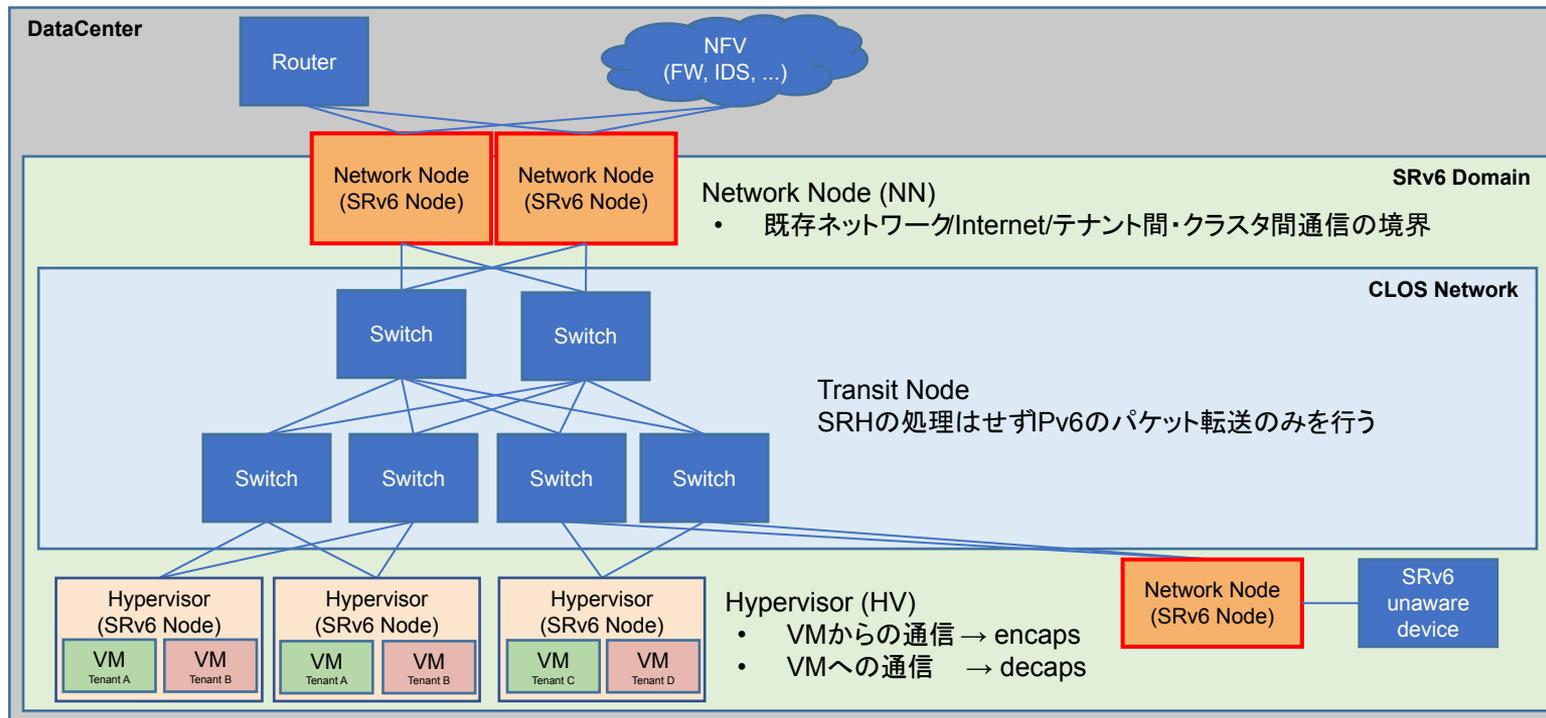
LINEにおける現在の構成



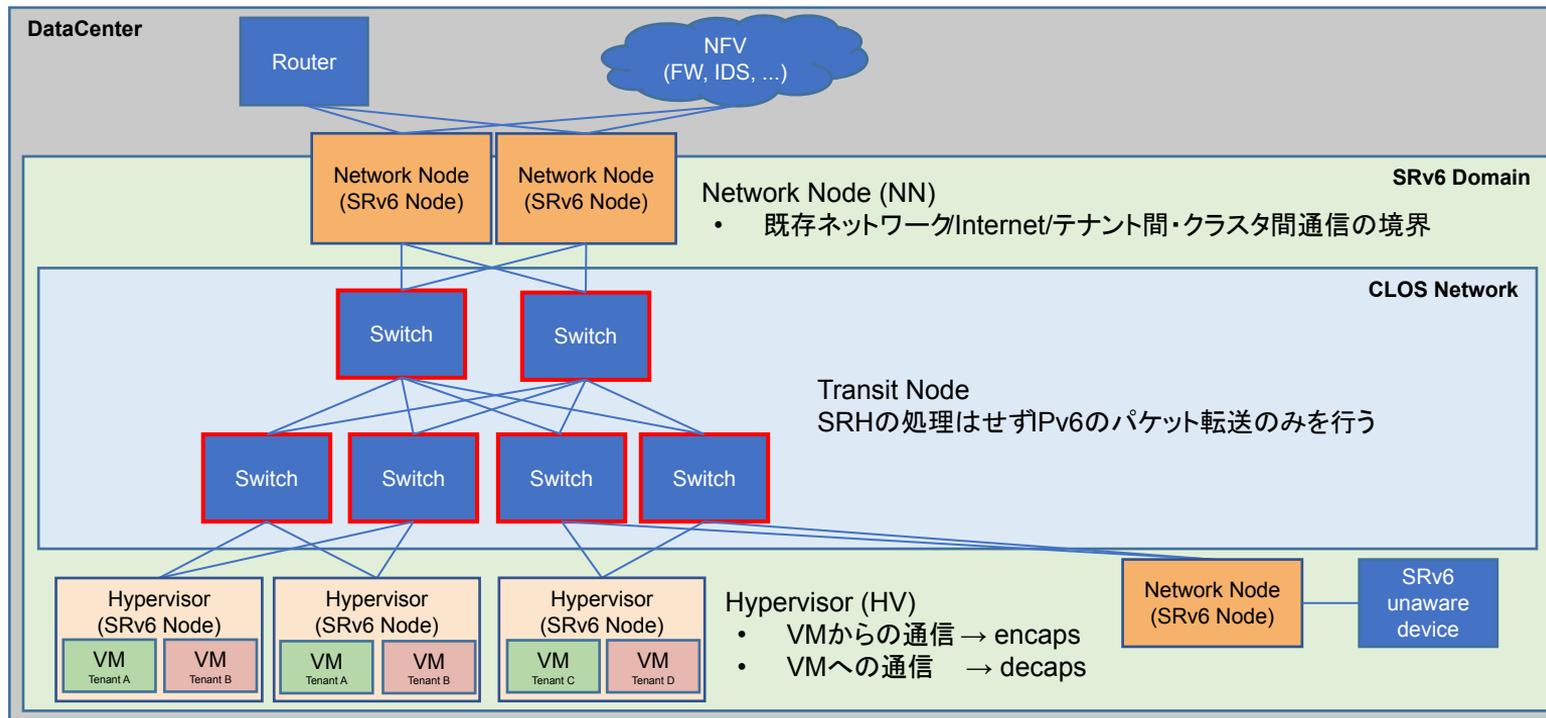
LINEにおける現在の構成



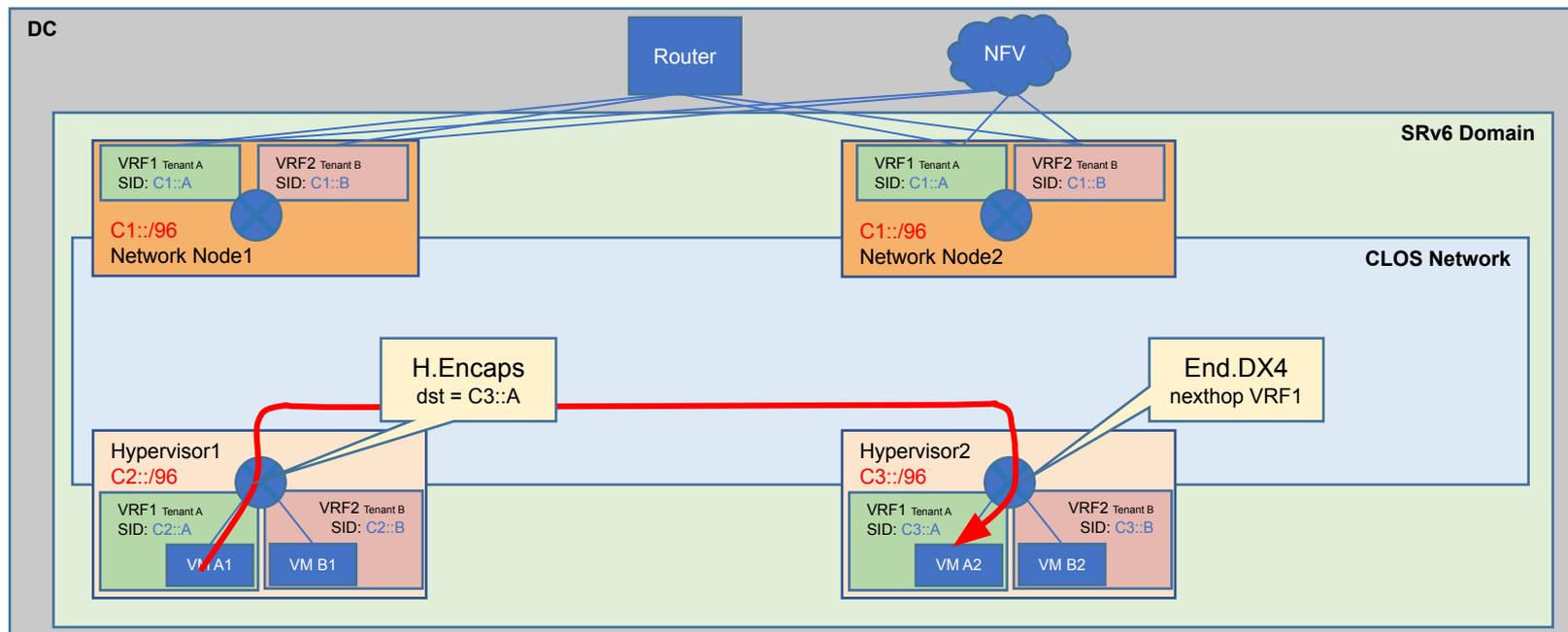
LINEにおける現在の構成



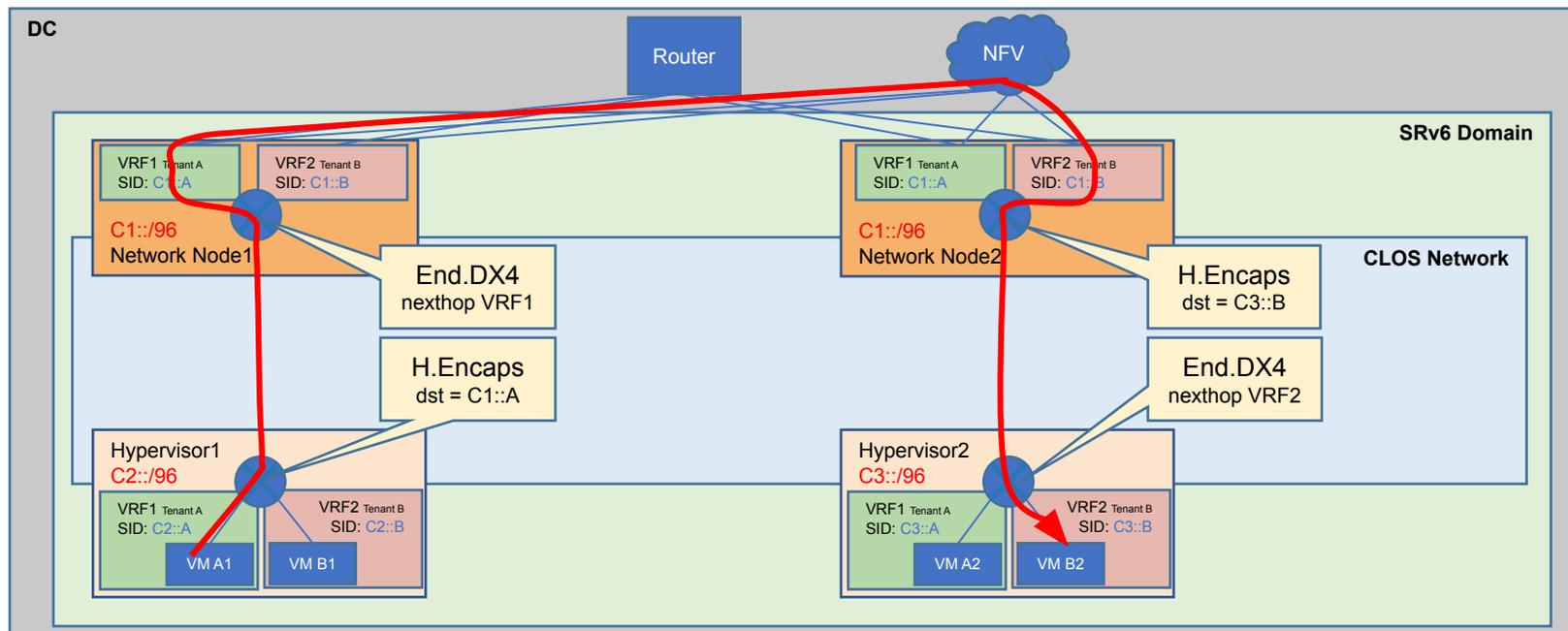
LINEにおける現在の構成



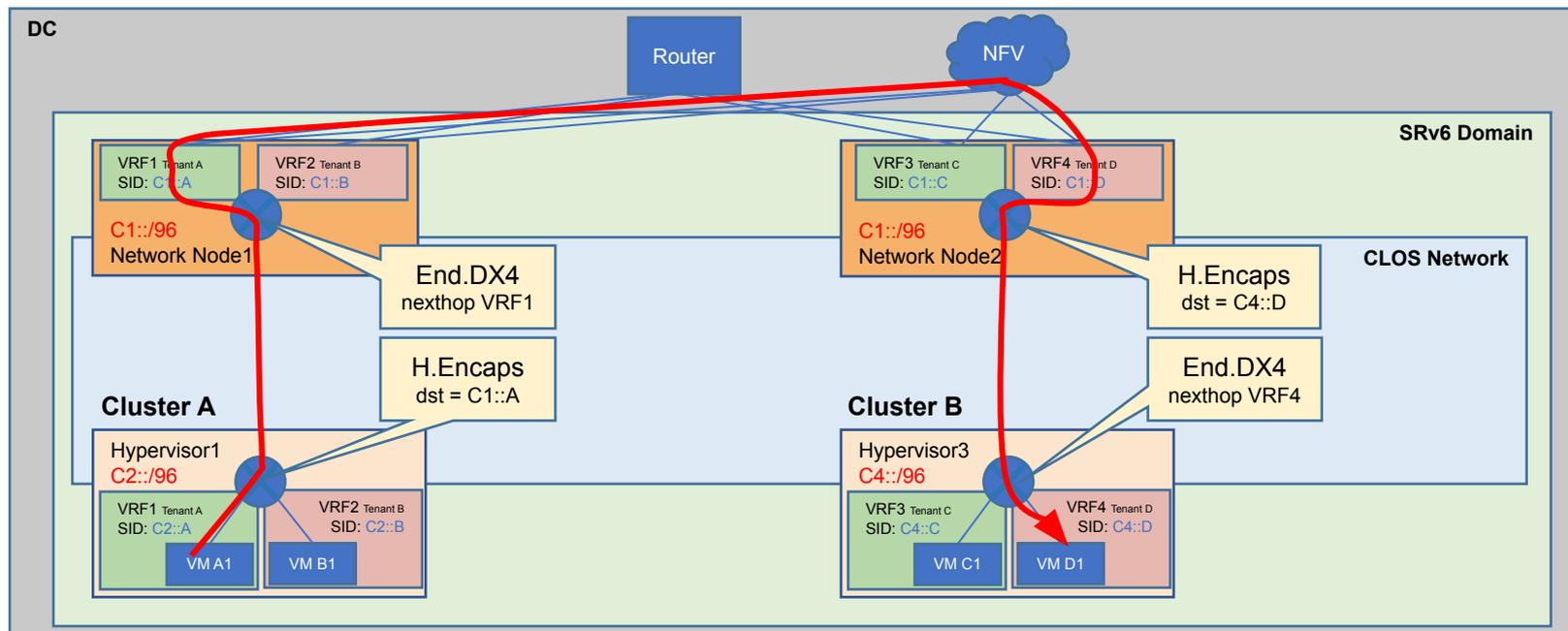
テナント内通信



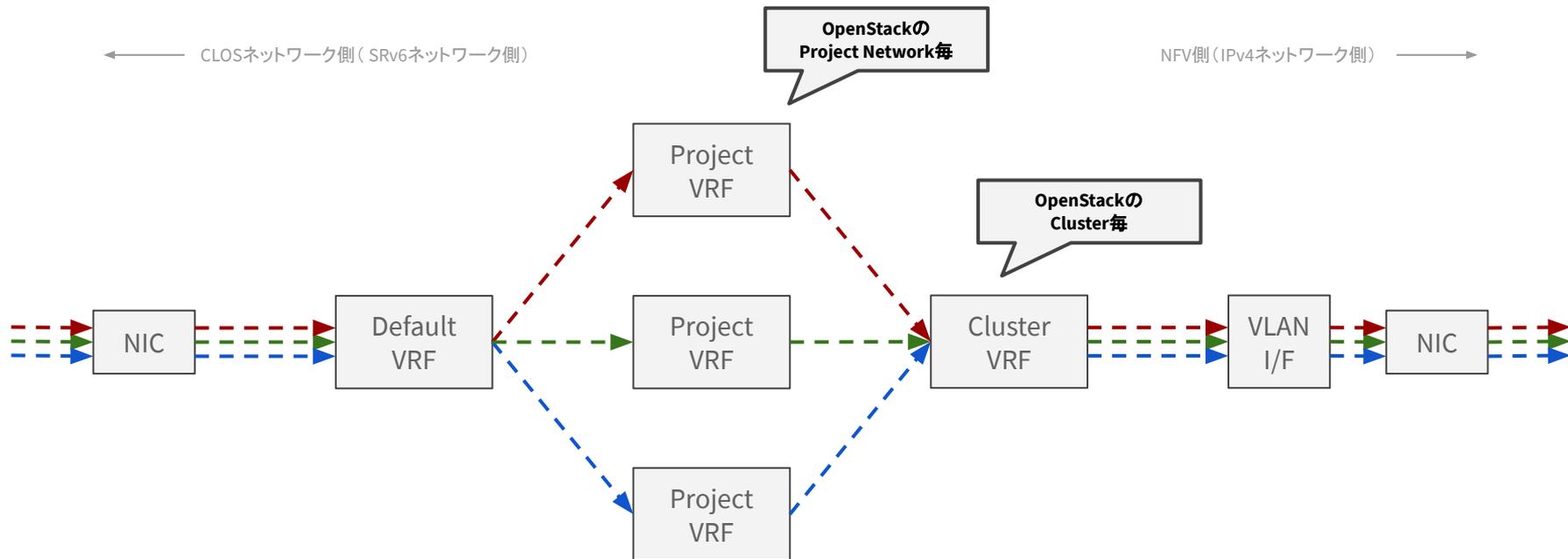
テナント間通信



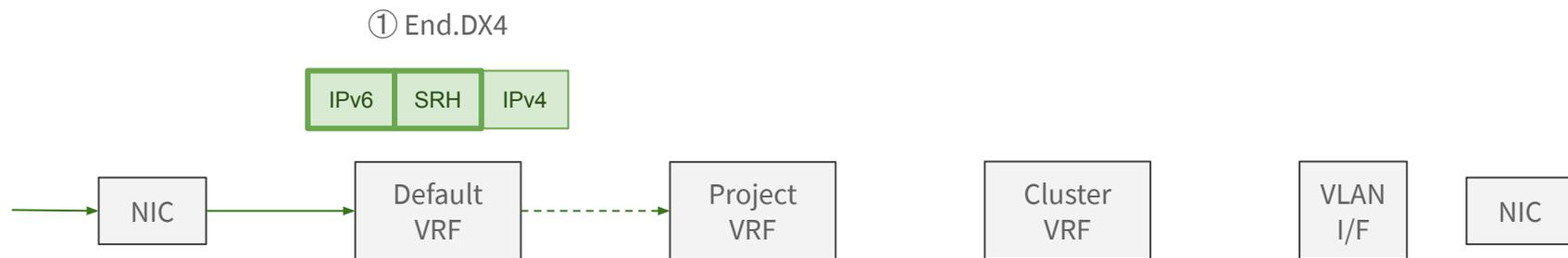
クラスタ間通信



ネットワークノードの処理

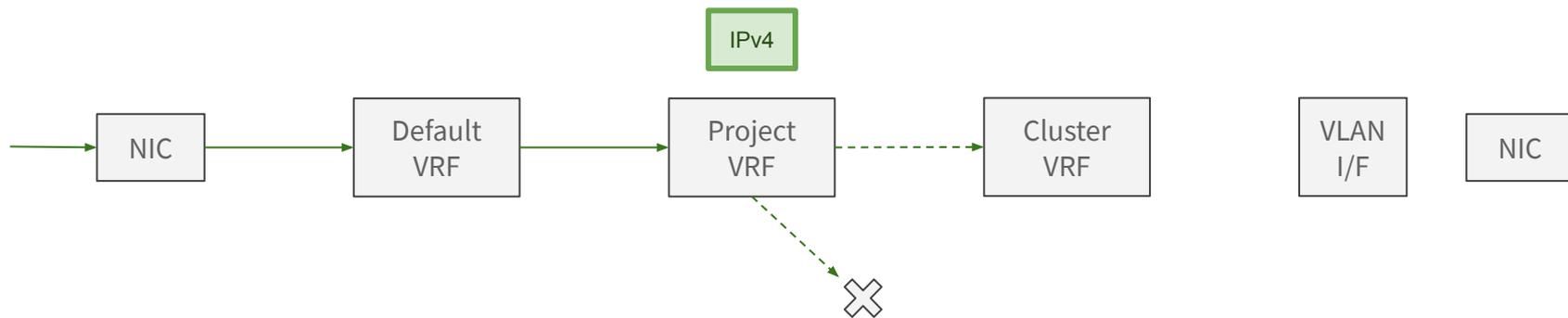


ネットワークノードの処理

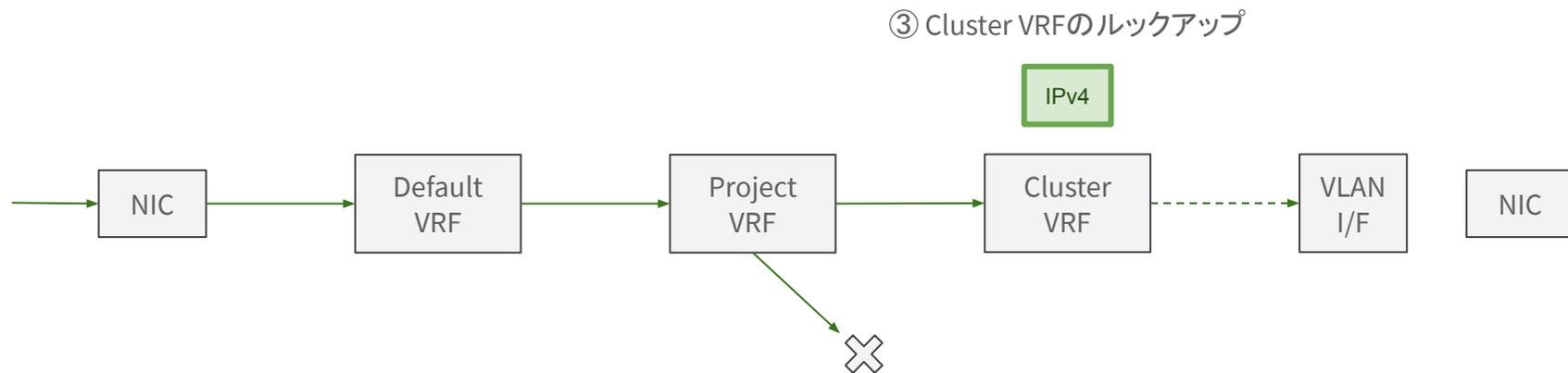


ネットワークノードの処理

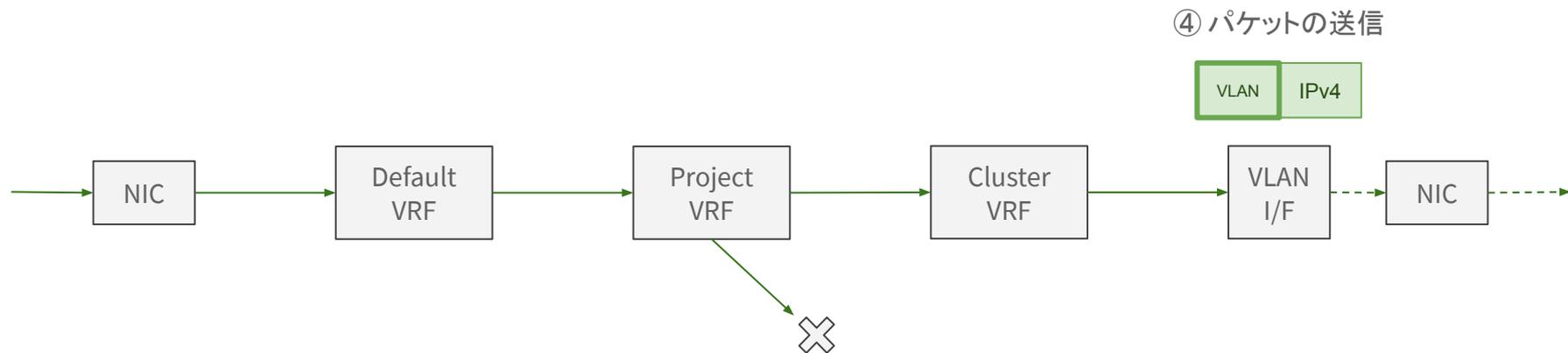
② Project VRFのルックアップ



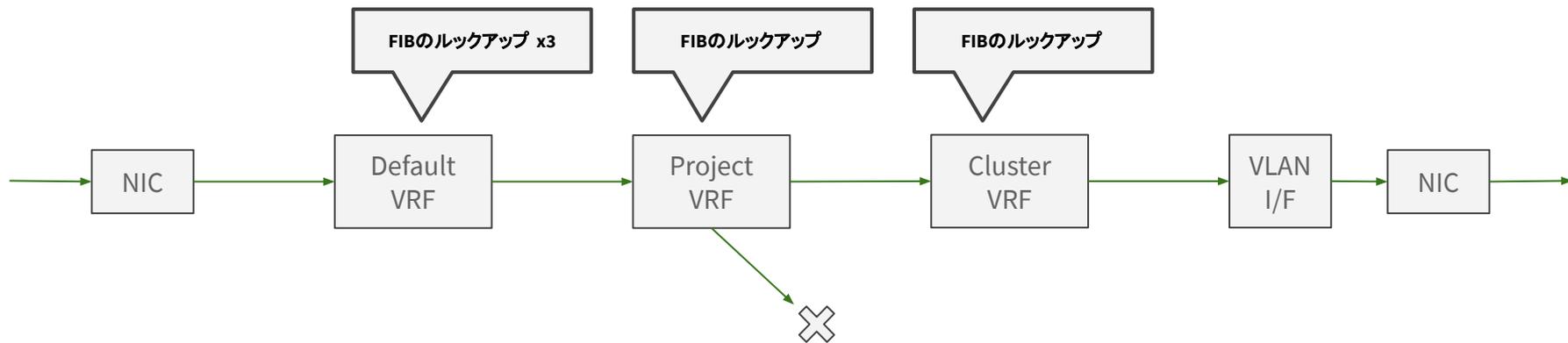
ネットワークノードの処理



ネットワークノードの処理



ネットワークノードの処理



問題点

- ネットワーク構成
 - LinuxのVRFの仕様やルーティングを駆使して頑張っているが、構成が複雑になる
- 性能問題
 - ソフトウェアでのFIBのルックアップは比較的重いコストになる
 - 何度もFIBのルックアップを行うため、普通の通信よりも性能が出ない
 - ネットワークノードを模した性能測定では普通のIPv6転送の半分程度しか性能が出ない
普通のIPv6転送: 1.09Mpps (Single Flow, 64bytes)
既存の構成: 0.51Mpps (Single Flow, 64bytes)

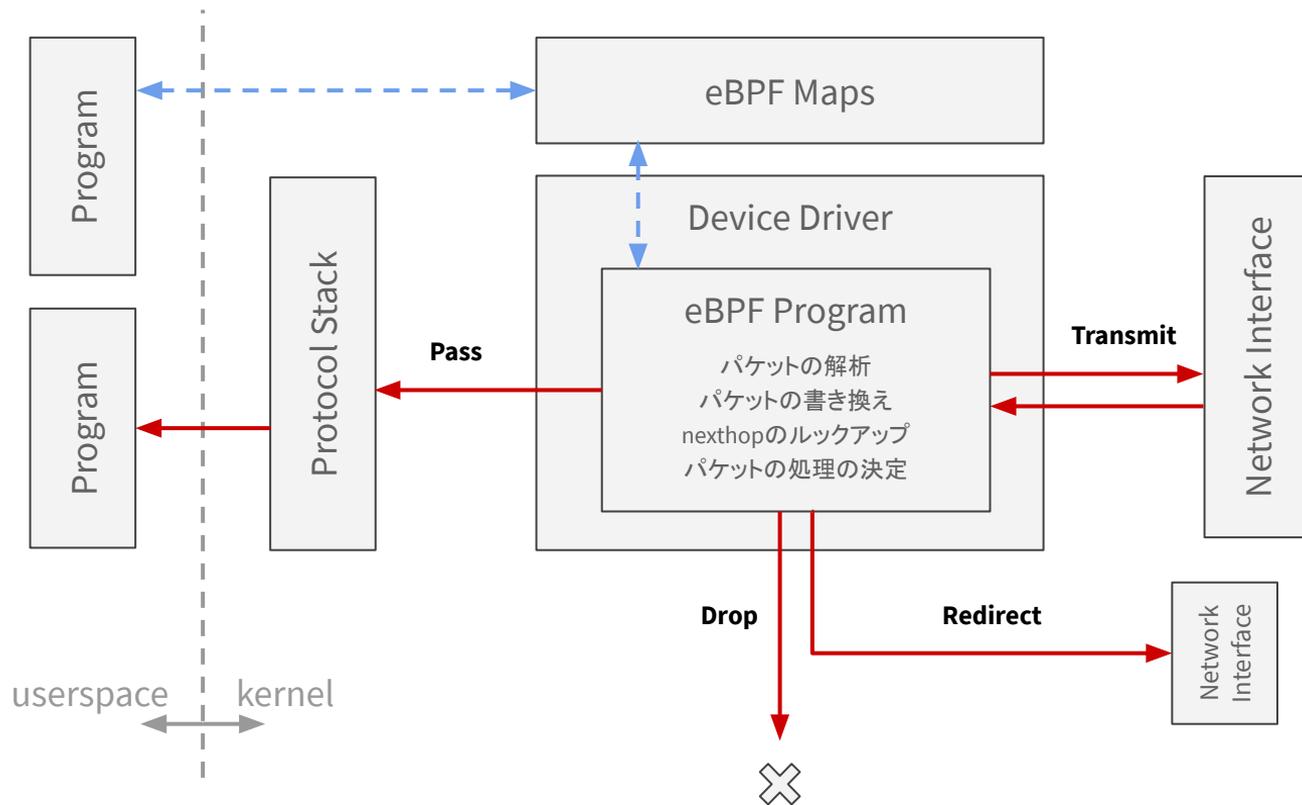
データプレーンの実装方法

- **XDP**
- DPDK (VPP)
- Hardware Switch
- ...



- L4LBでの使用実績がある
- ARPやICMPの処理をカーネルに依頼することができる
- LinuxカーネルのFIBを利用することができ、既存の構成からの変更点を少なくすることができる

XDP (eXpress Data Path) の概要

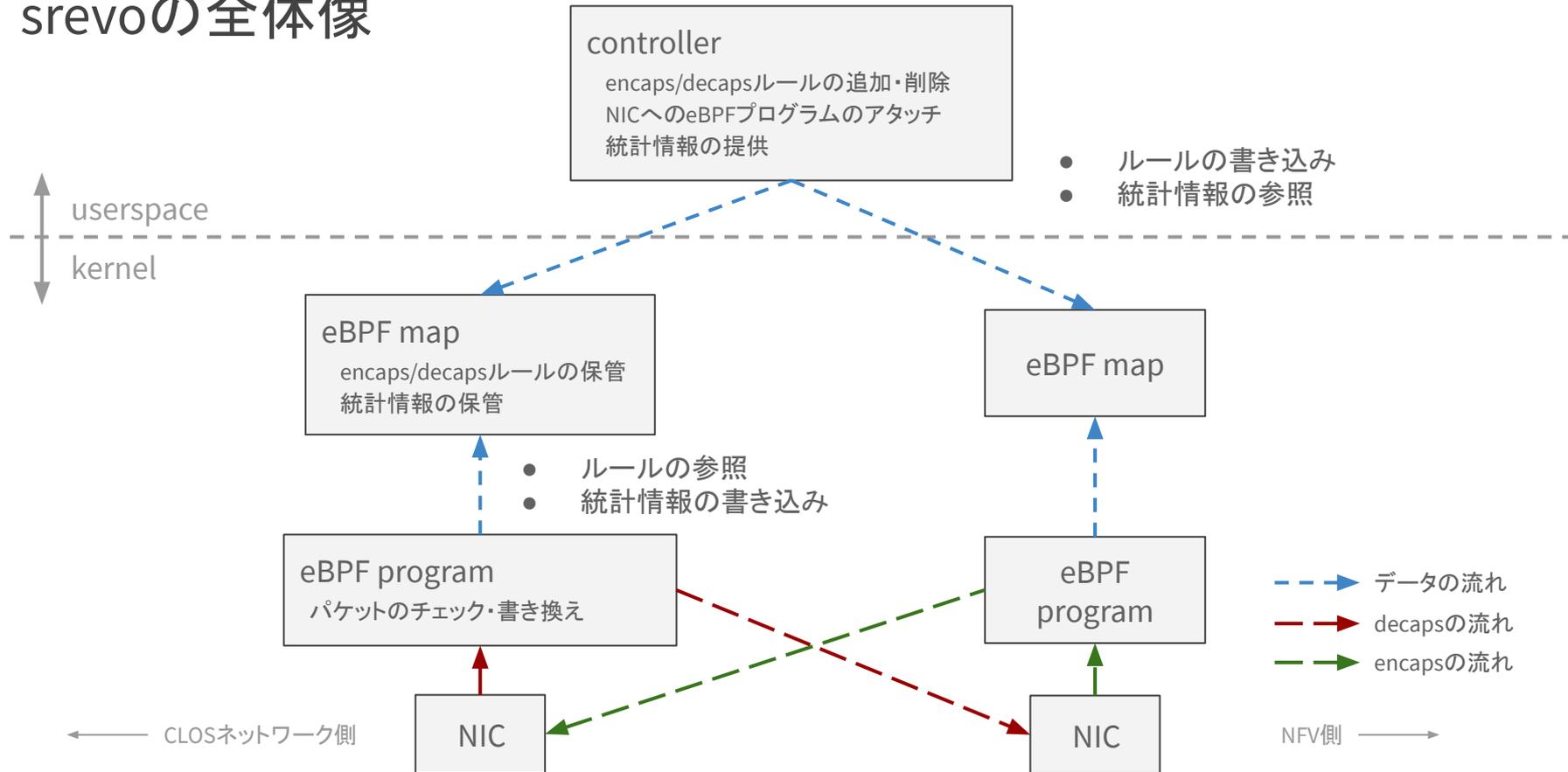


srevo : Faster SRv6 D-plane with XDP

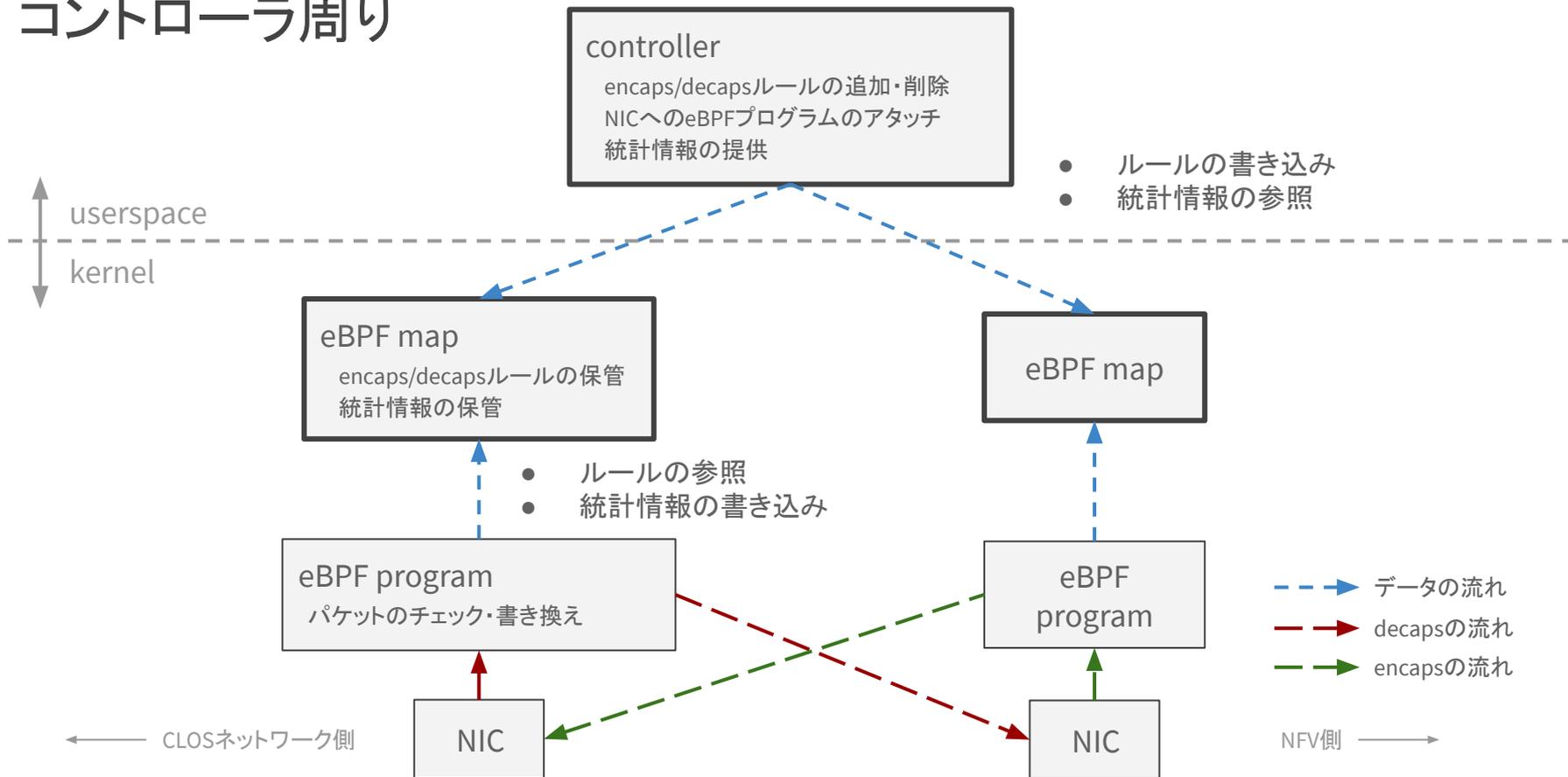
srevoの概要

- XDPを用いたSRv6のパケット処理基盤
 - さまざまな最適化を行い、既存の方法よりも高速に SRv6のパケットを処理できる
- LINEの用途に特化させることでシンプルな実装
 - 必要のない機能は実装を省いている
 - データプレーン自体は約 500行

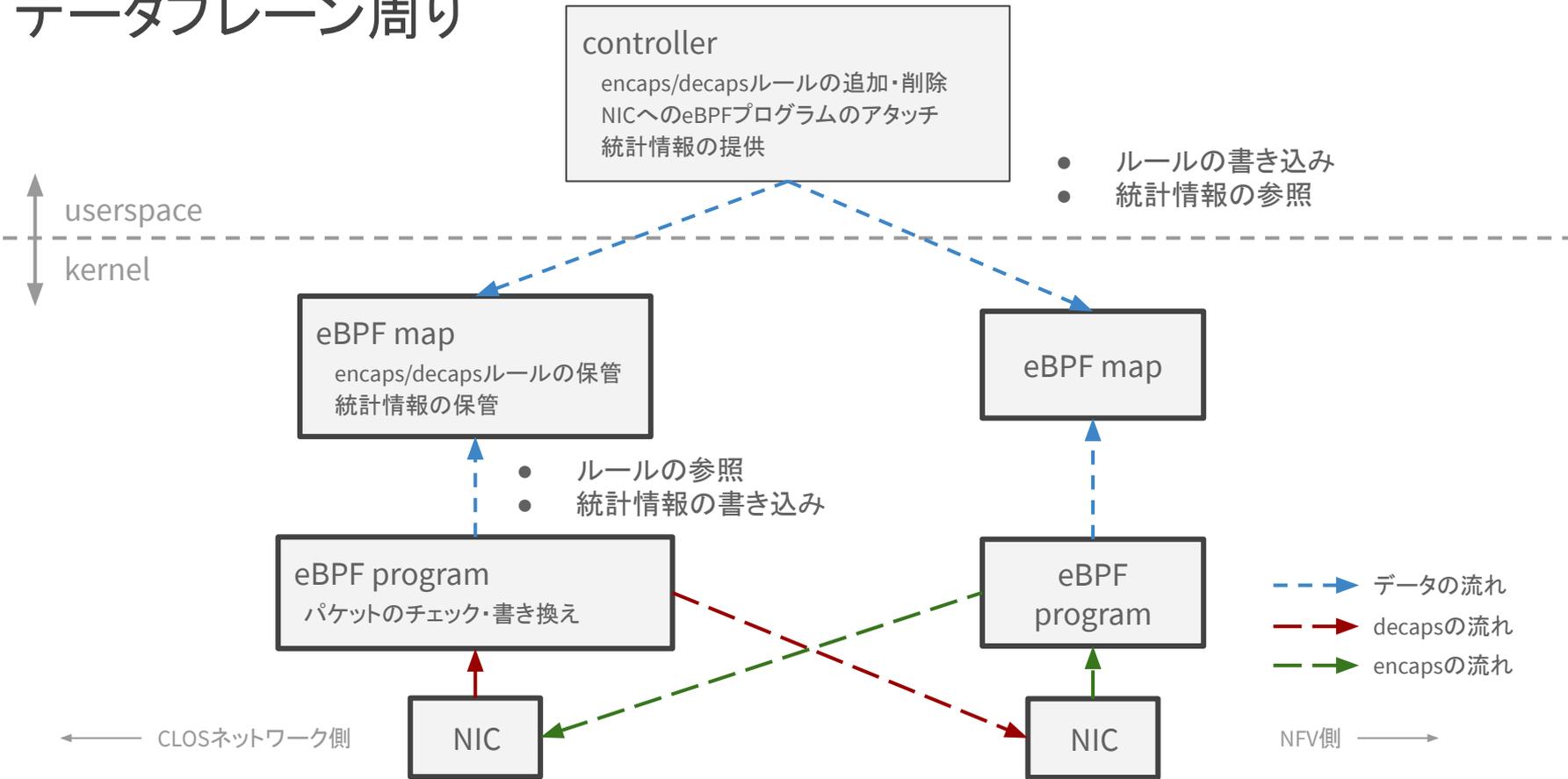
srevoの全体像



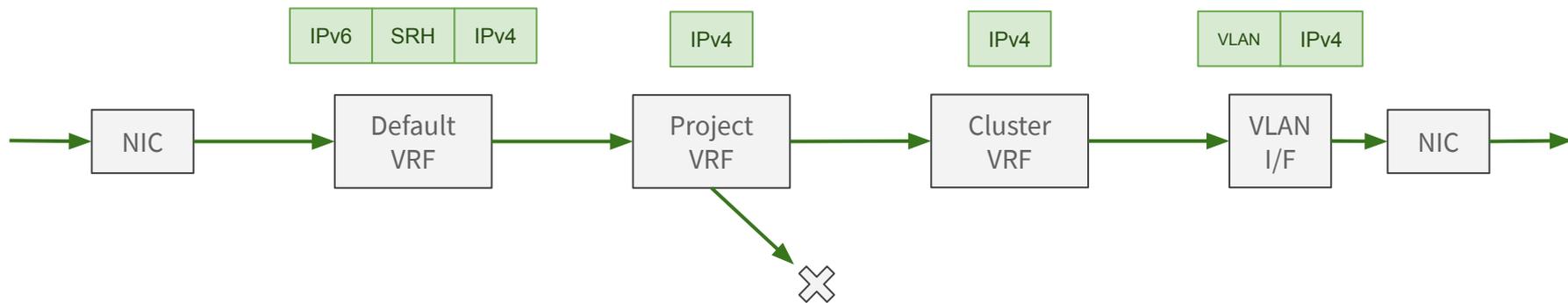
コントローラ周り



データプレーン周り



最適化の方針

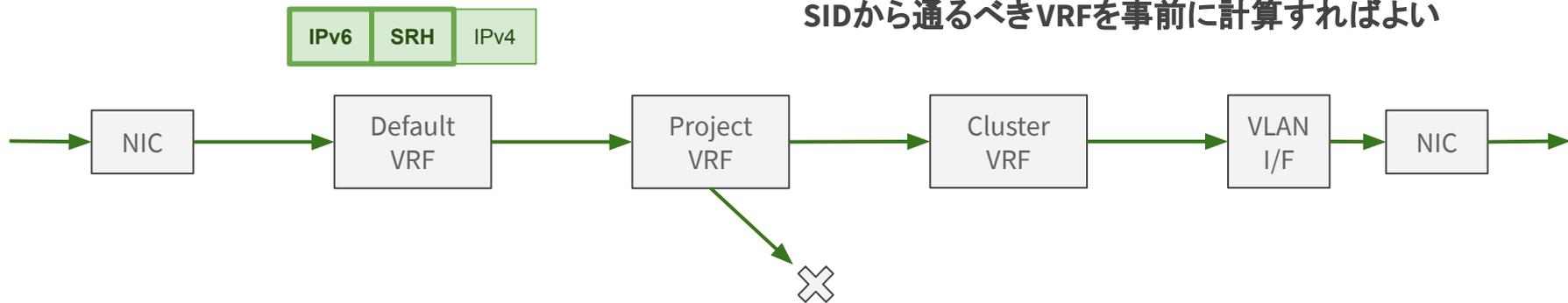


最適化の方針

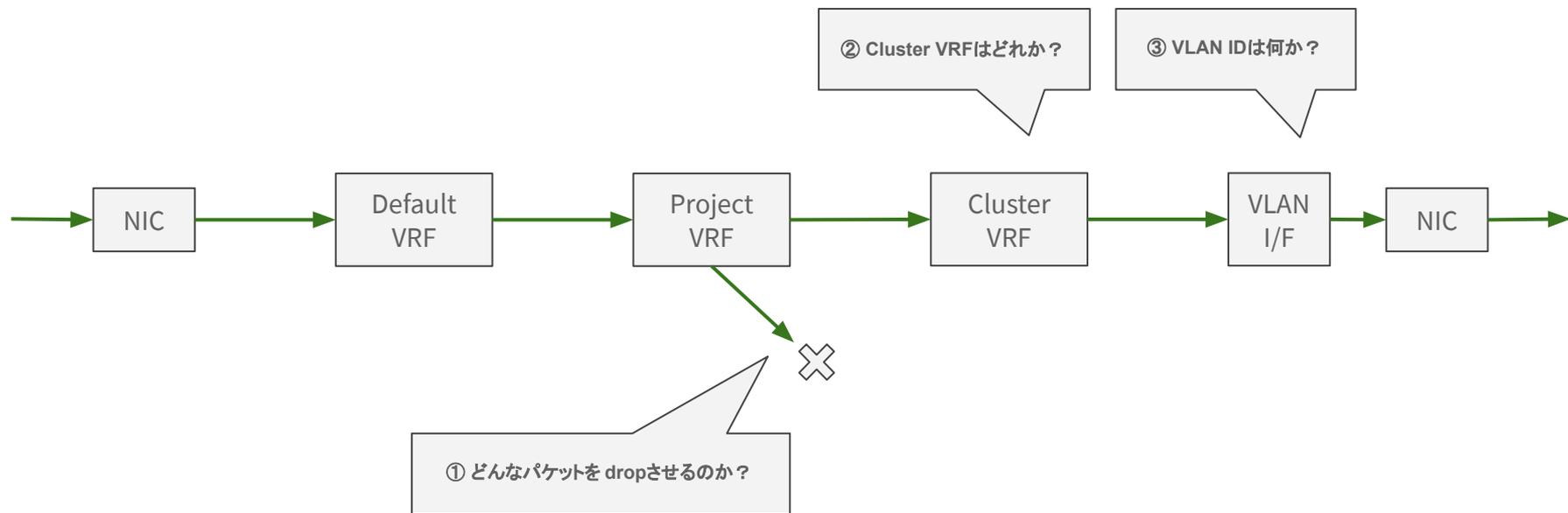
SIDが決まれば通るべきVRFが決定する
SIDとVRFの関係は一度決まれば変更されない



SIDから通るべきVRFを事前に計算すればよい



最適化の方針



データプレーンの処理

① 処理すべきかチェック

eBPF map



データプレーンの処理

② eBPF mapのルックアップ

eBPF map

SID	Cluster VRF	VLAN I/F
2001:db8::1	A	A
...

IPv6 SRH IPv4

NIC

Default VRF

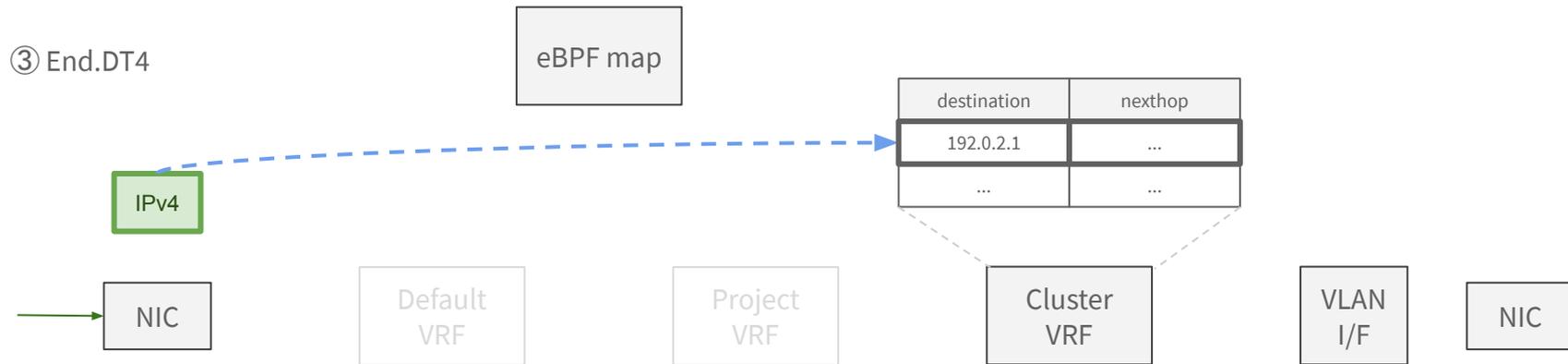
Project VRF

Cluster VRF

VLAN I/F

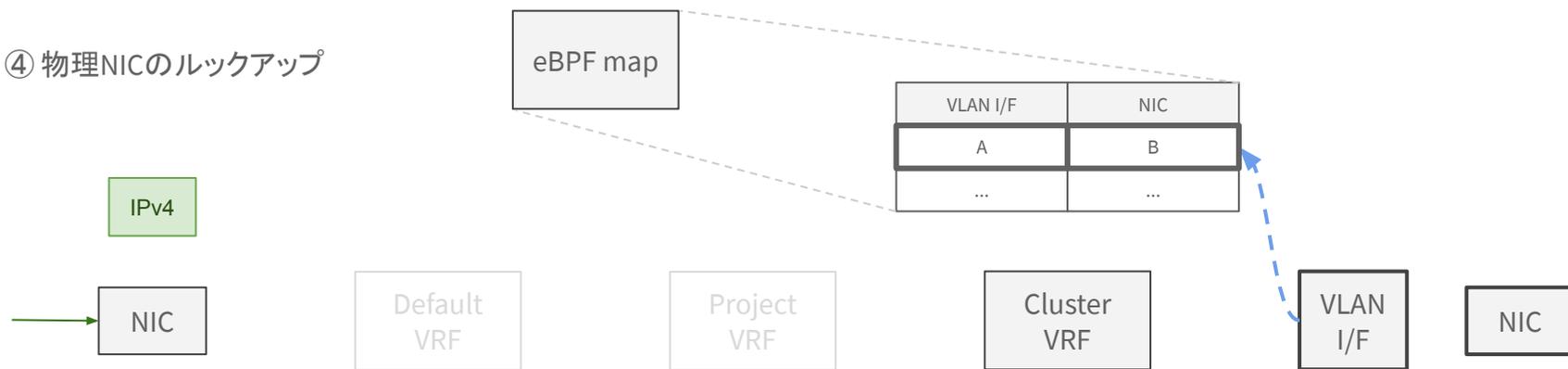
NIC

データプレーンの処理



データプレーンの処理

④ 物理NICのルックアップ



データプレーンの処理

⑤ VLANヘッダの挿入

eBPF map

VLAN IPv4

NIC

Default VRF

Project VRF

Cluster VRF

VLAN I/F

NIC

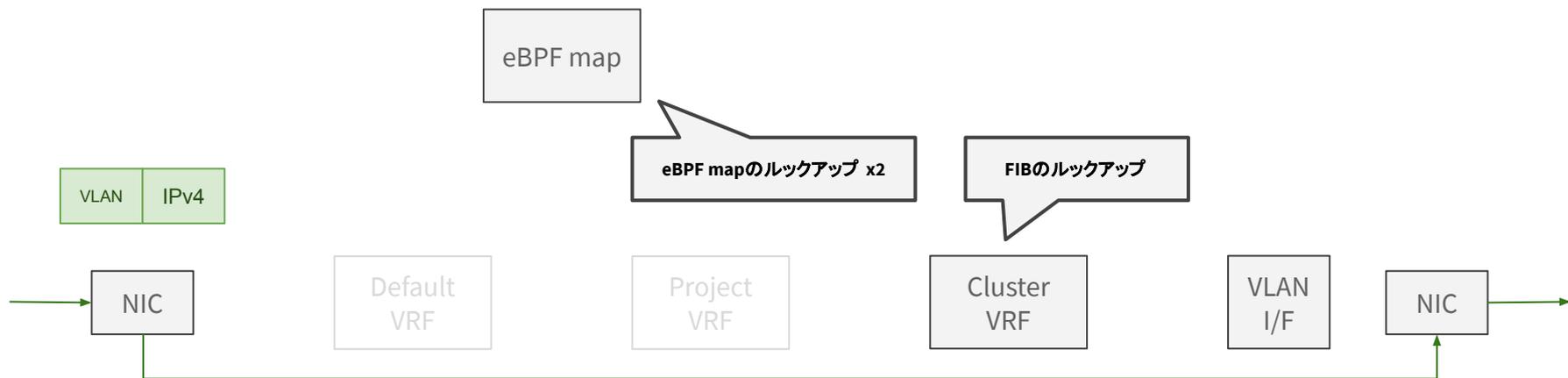
データプレーンの処理

⑥ XDP_REDIRECT

eBPF map



データプレーンの処理



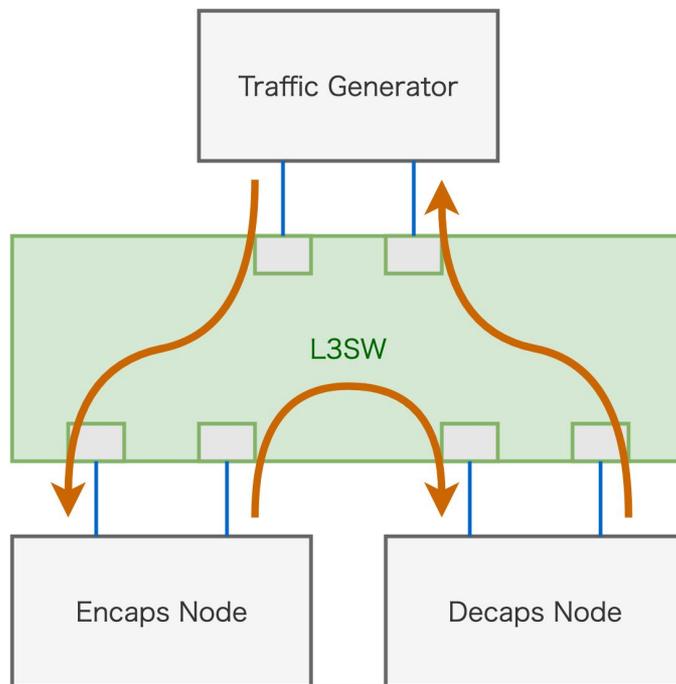
最適化の内容

- FIBのルックアップを減らすことが出来た
 - 5回 => 1回に減らすことができた
 - 代わりにeBPF mapを2回ルックアップする必要がある
- Project VRFを減らすことが出来た
 - VRFで頑張ってたEnd.DT4していた部分が無くなったので余計な VRFが要らなくなった

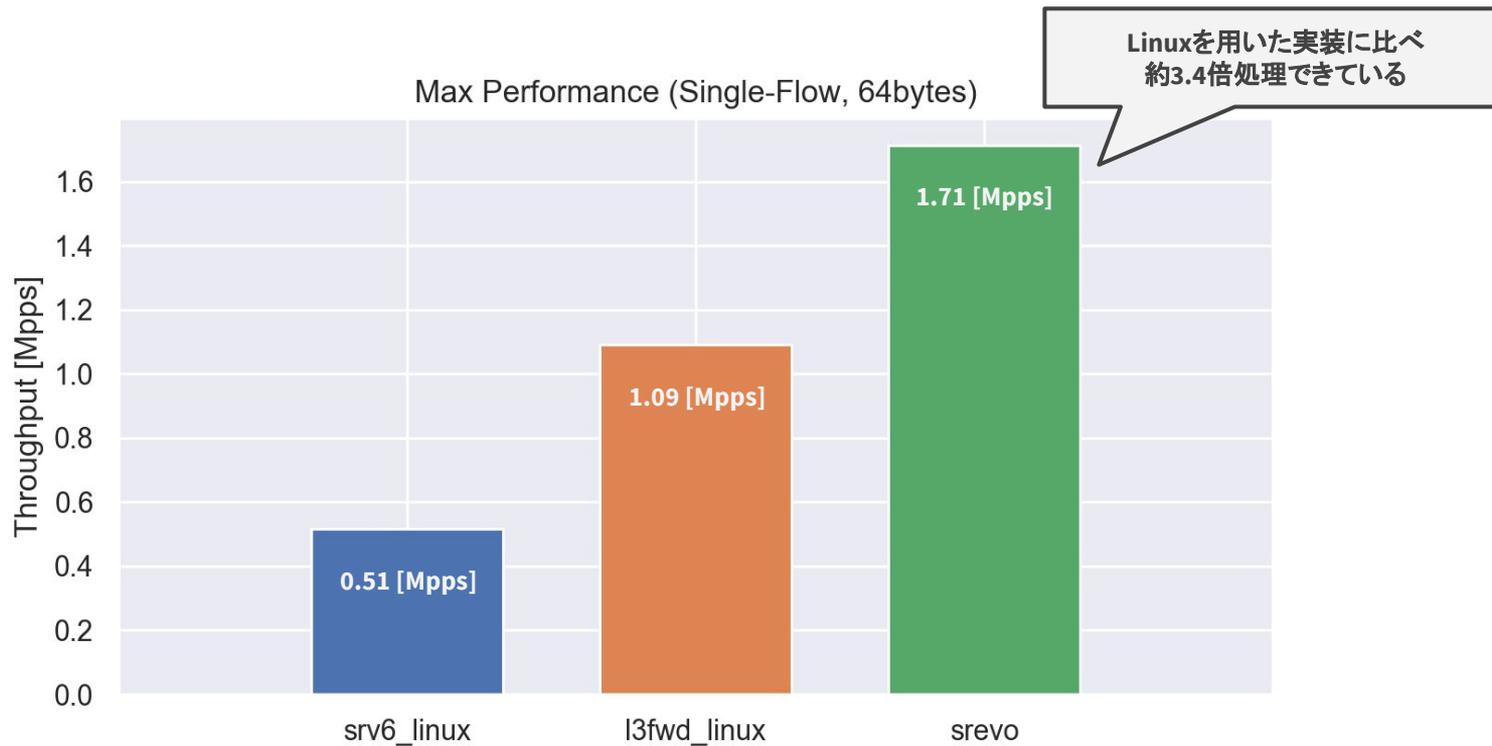
Performance

測定環境

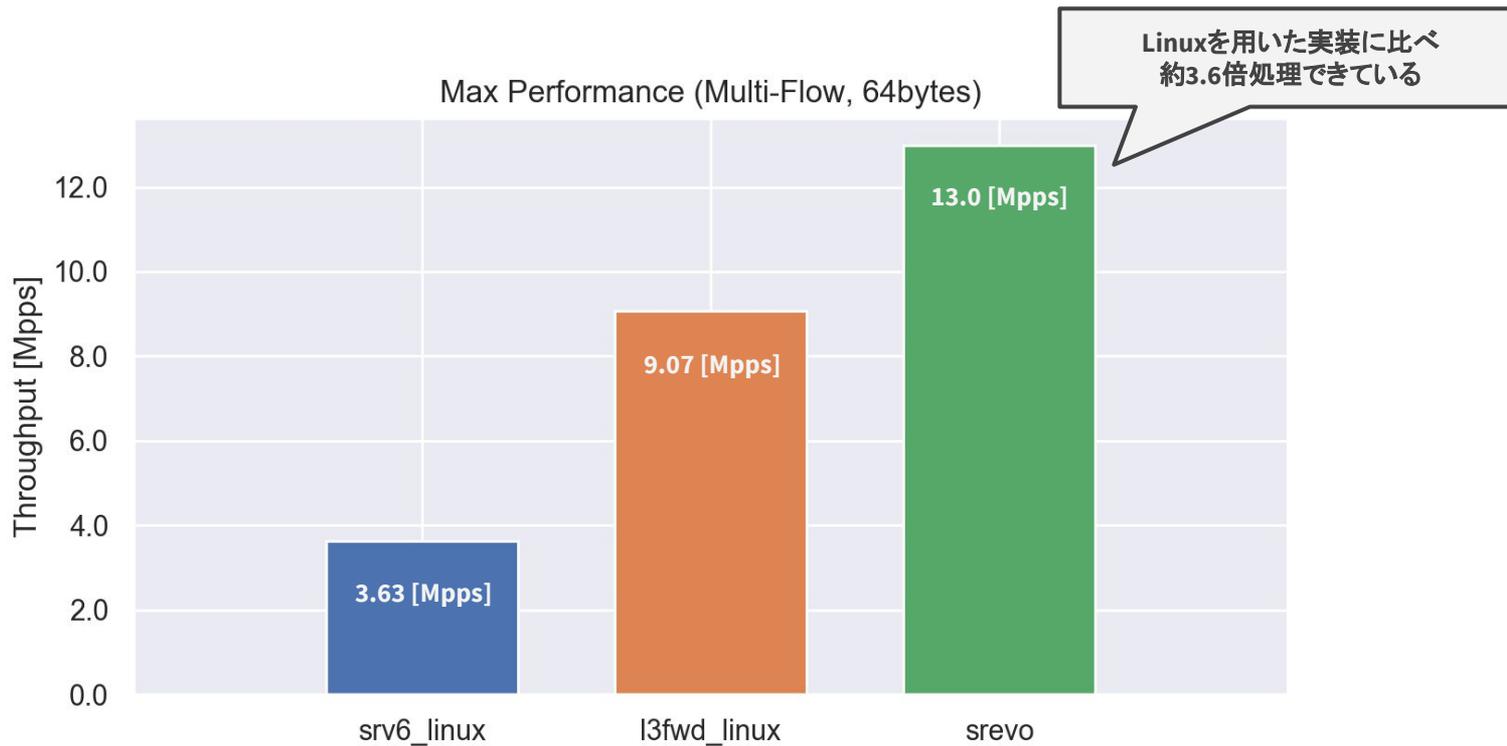
- Traffic Generator
 - Intel Xeon E5-2630 x2, 256GB
 - Intel XL710 40GbE 2port
 - Linux 3.10
 - TRex v2.50
- Encaps Node / Decaps Node
 - Intel Xeon Silver 4114 x2, 128GB
 - Mellanox MCX516A-CCAT 100GbE 2port
 - Linux 5.3 rc5
- L3SW
 - Mellanox SN2700
 - Cumulus Linux 3.5.3



測定結果



測定結果



XDPを選択してよかったこと

- Linuxカーネルの恩恵を受けることができる
 - ARPの処理やFIBの管理はカーネルに依頼することができる
- 必要な機能だけ実装するのでシンプルに保つ
 - データプレーンのプログラム自体は 500行程度
 - 新しい機能が欲くなれば実装すればよくなる
- 既存のコントロールプレーンを利用できる
 - 大規模な変更を加えることなく、今動いている SRv6の仕組みにそのまま乗られる
 - srevoが問題を起こした際のフォールバックとして既存の実装がそのまま利用できる

苦労した点

- デバッグがしづらい
 - tcpdump等のデバッグツールが使えないのでパケットが消えたときにすごく困る
- カーネルのすべての情報に触れるわけではない
 - FIBのルックアップをしてもVLAN I/Fの情報しか取れなく、物理 NICとの対応関係は自分で管理する必要がある

今後の課題

- デバッグブルなプログラムにする
 - 監視用のインターフェースの実装
 - デバッグモードの実装
- Service Function Chainingに対応する
 - 複数のSIDでencapsできるようにする
 - End.DT4以外のFunctionの実装 (EndやEnd.Xなど)
- さらなる最適化
 - decaps / encapsの処理の高速化
 - ルールを格納するデータ構造の持ち方

議論したいポイント

- この取り組み(DCNにSRv6を用いることやXDPを用いた最適化)についてどう思いましたか？
- SRv6を用いてどのような世界が実現できると嬉しいと思いますか？
- データプレーンのコントローラにどんな機能があると嬉しいと思いますか？(デバッグ機能など)