

# 今こそ学びたい Kubernetesネットワーク

## ～CNIが繋ぐNWとプラットフォームの「フラット」な対話



Takuto Nagami  
@logica0419




Kotaro Kawasaki  
@n4mlz

# 自己紹介

- Takuto Nagami (@logica0419)
- 千葉工業大学 情報科学部 情報ネットワーク学科 4年
  - 研究: **合意アルゴリズム** (Raft) × **暗号** (秘密分散)
- 得意技術: Go言語 / コンテナ / Kubernetes
- JANOG歴
  - JANOG56 **若者支援**で参加
  - JANOG57 初登壇！



# 自己紹介

- Kotaro Kawasaki (@n4mlz)
- 筑波大学 情報学群情報科学類 3年
- 趣味: コンテナランタイム、OS 自作、車輪の再発明
- 最近  **Cyrius** という Linux コンテナをネイティブ実行できる自作 OS を Rust で書いています



モダンなインフラ  
チェックしていますか？



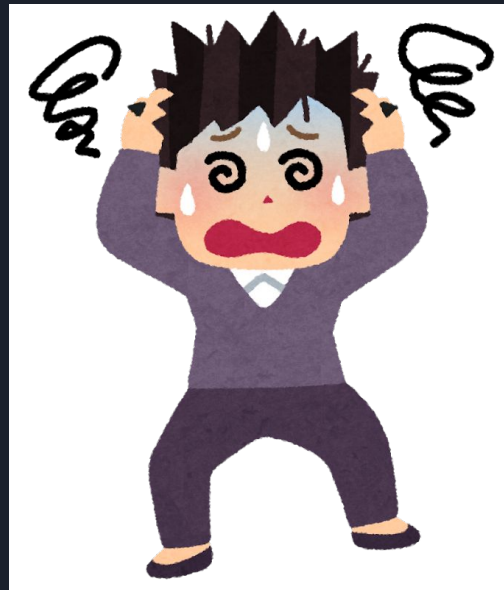
# クラウドからクラウドネイティブへ

- ここ十数年で、**クラウド**ビジネスは大躍進を遂げた
- クラウドの設計思想は**クラウドネイティブ**へと昇華
  - クラウドで一般的になった、**スケールしやすい**開発アプローチを**あらゆる場所**で活用する
  - **オンプレ** (プライベートクラウド) でも、この思想が導入されつつある
- このクラウドネイティブなエコシステムの中心にるのが**Kubernetes (K8s)**



# K8s、正直わかりませんよね

- 「**難しい**」というイメージが先行している
- やはり「プラットフォームは**我々の管轄外**」という考えをしている人が多いのでは



# プラットフォームチームも気持ちは同じ

- プラットフォーム「ネットワーク？**守備範囲**じゃない」
- ネットワーク「K8s？**守備範囲**じゃない」





## ところで: 架け橋は用意されている

- K8sのネットワークは**差し替え可能**な設計
  - K8sのオリジナル開発者はGoogle
  - 自身のクラウドのネットワークに**最適化**できるようにあえて**どんなネットワーク技術も取り込める**ように設計している
- プラットフォームとネットワークを繋ぐ**開かれた門**かつ**共通言語**の役割を担う部品がいくつかある



# しかしそれらの部品の実態は…

- ネットワークもプラットフォームもお互いが忌避し、  
両者の**断絶を生む閉じた門**になっている





しかしそれらの部品の実態は…

**なんて  
もったいない！**





# OSSも銀の弾丸ではない

- **オンプレ**であっても、**OSS**を組み合わせでK8sのネットワークを構成している現場がほとんど
  - ネットワークチームは**アンダーレイのみ**管理し、  
K8sのネットワークは**別チーム**がOSSで運用
- OSSは**汎用性**のため、**SDN**ベース
- しかし、SDNより**物理機材**ベースの方が**効率が良い**はず
  - OSSに頼り切るのではなく、**ネットワークチーム**が  
**介入して自作・改造**することで効率化が図れる



# K8sネットワークの自社開発・改造例

- **Google Cloud:** GKE Dataplane V1 / V2
- **AWS:** Amazon VPC CNI
- **Microsoft Azure:** Azure CNI
- **Cybozu:** Coil on Neko基盤
  - 国内のクラウドでの例はこれしか見たことがない  
(公開されていないだけかもしれません)
  - ↑ この状況を変えよう！というのが今回の目的




# プラットフォームとフラットと対話する

- ネットワークとプラットフォームの対話なしに、真に  
効率の良いネットワークはあり得ない
- 今回のテーマ「フラットとJANOG」に合わせ、フラットと  
気兼ねない対話に必要な**共通言語**を身につけましょう！



**JANOG57**  
in OSAKA



# 「プラットフォームと ネットワークの フラットな対話」を作る

これを今回のゴールとします！一緒に頑張りましょう



# アジェンダ

- イン트로ダクション ←イマココ
- Kubernetesを知る
- Kubernetesとネットワーク
- Linuxのネットワークスタック
- L2/L3ネットワーク: CNI
- L4ロードバランサー: kube-proxy
- プラットフォームとの対話で広がる可能性
- まとめ

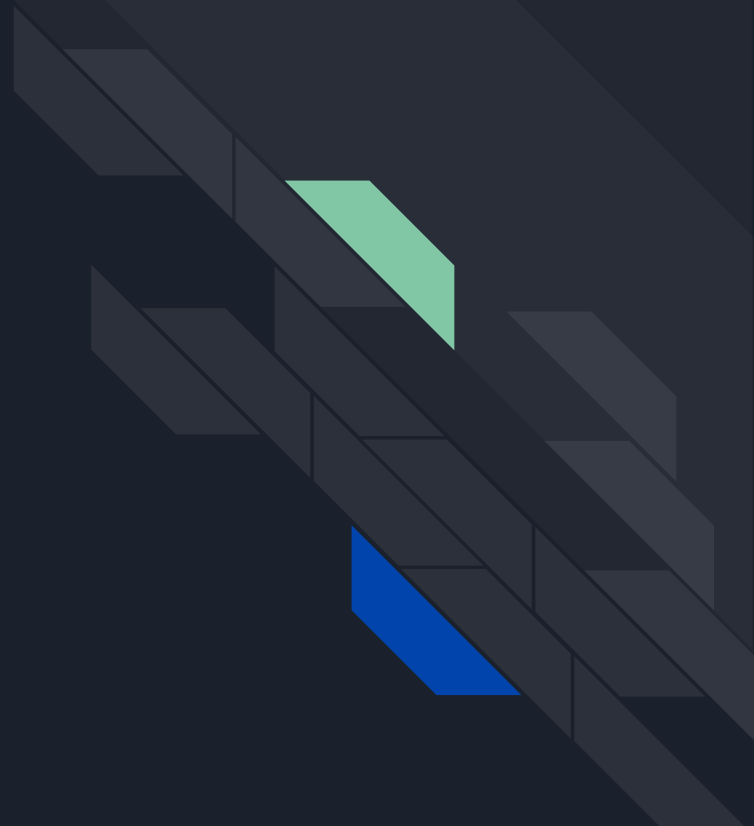



# アジェンダ

- イントロダクション
- Kubernetesを知る ←イマココ
- Kubernetesとネットワーク
- Linuxのネットワークスタック
- L2/L3ネットワーク: CNI
- L4ロードバランサー: kube-proxy
- プラットフォームとの対話で広がる可能性
- まとめ



# Kubernetesを知る

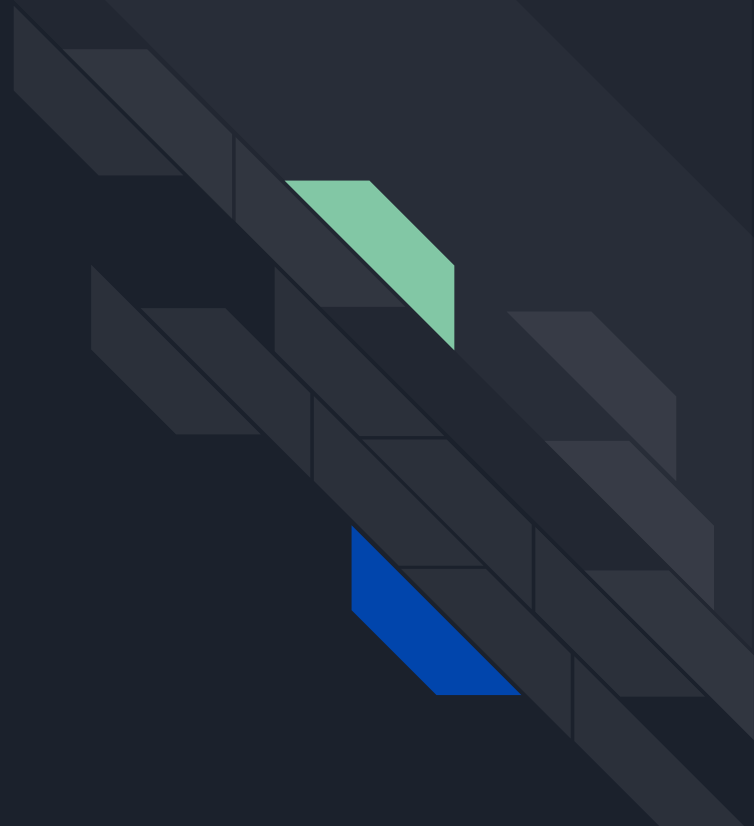




K8sは難しく思われがちだが  
**ポイントを絞れば**  
**そこまで難しくない！**

ひとまずKubernetesの世界を旅してみましょう

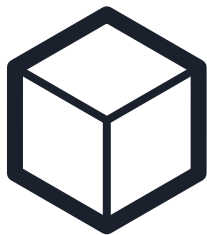
# Kubernetesの 提供するリソース



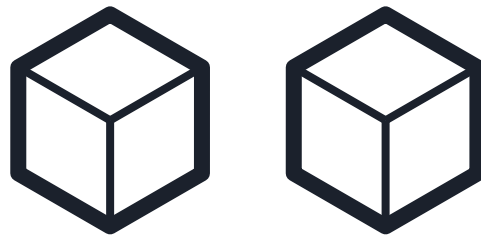
# Pod

- 便利になったコンテナ、**K8sの最小単位**
- 1つのワークロードを構成する単位
  - ≡ **1つのIPアドレス**でくくって良い範囲

Pod

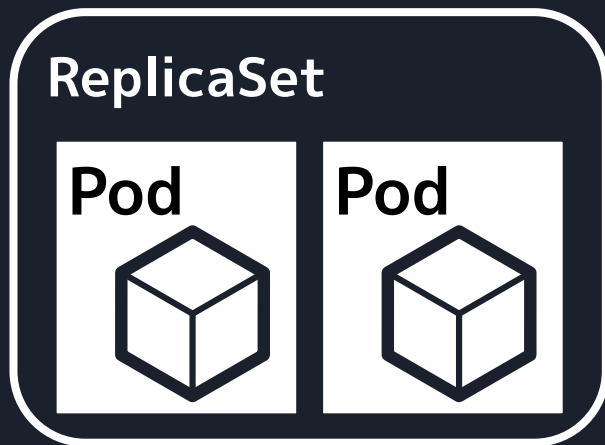


Pod



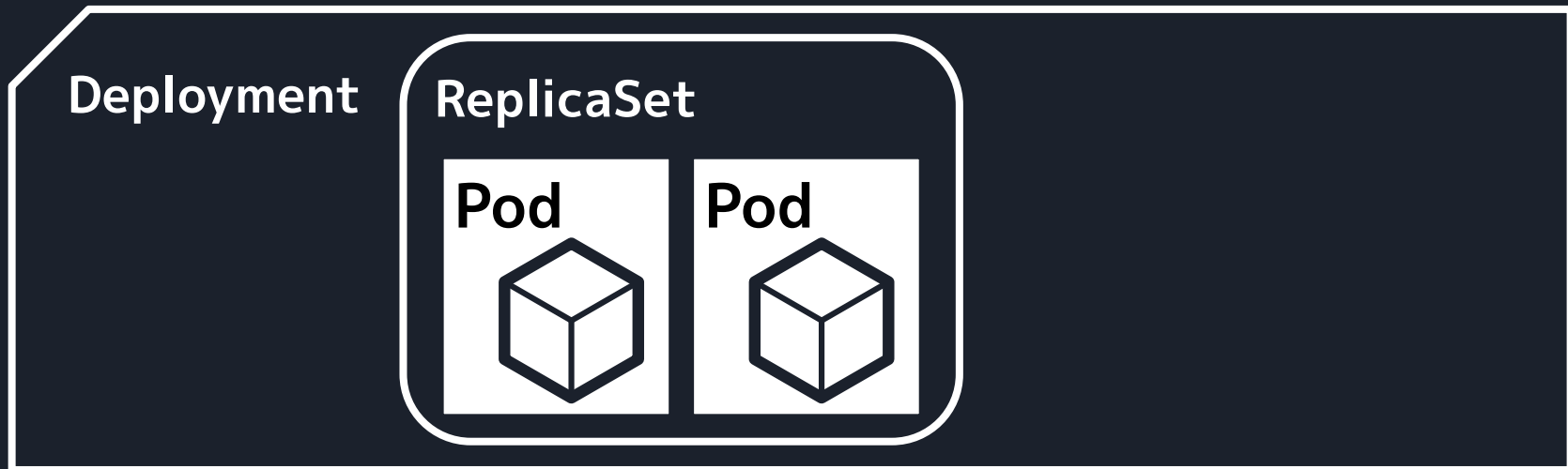
# ReplicaSet / Deployment

- Podをまとめて管理するための仕組み
- **ReplicaSet**: 同じ内容のPodを指定個数立てる
- **Deployment**: ReplicaSetを使い安全にバージョン移行



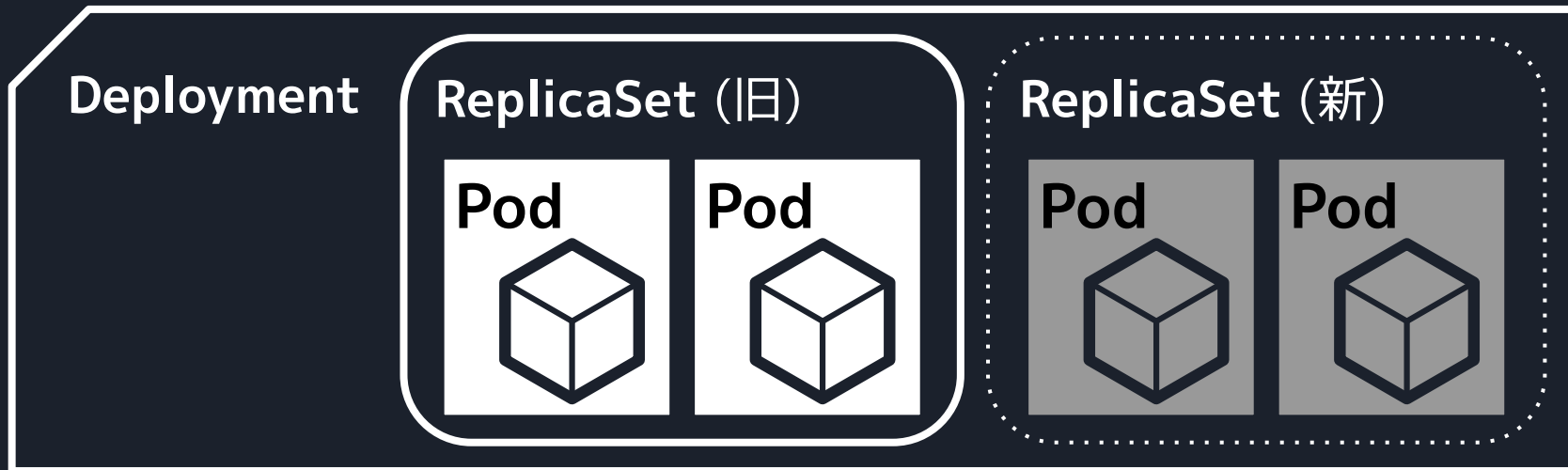
# ReplicaSet / Deployment

- Podをまとめて管理するための仕組み
- **ReplicaSet**: 同じ内容のPodを指定個数立てる
- **Deployment**: ReplicaSetを使い安全にバージョン移行



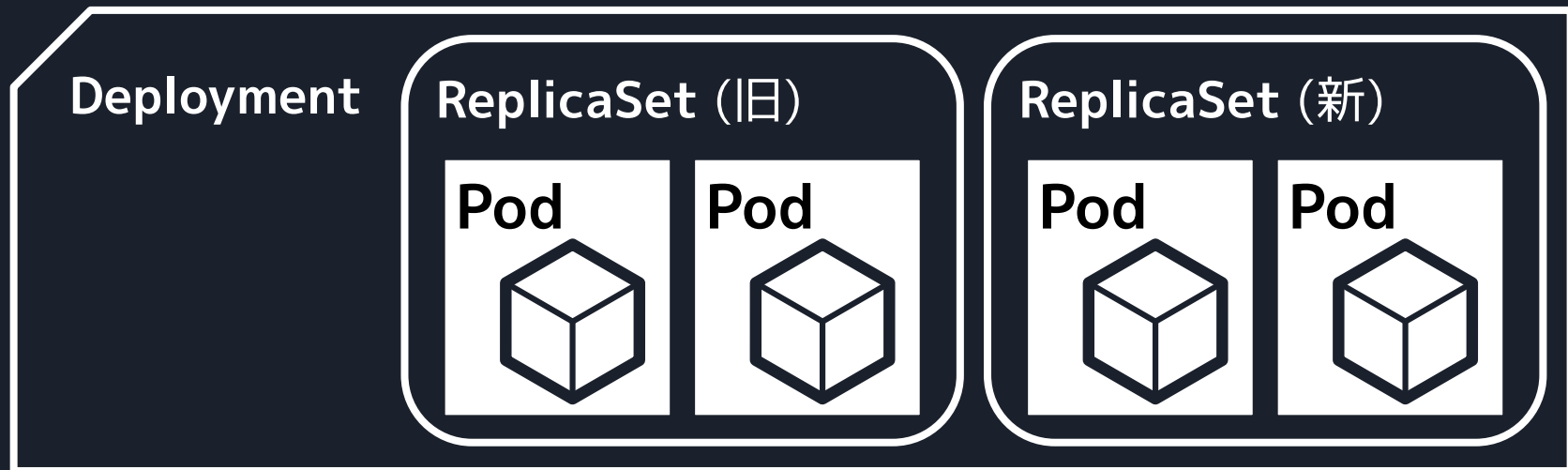
# ReplicaSet / Deployment

- Podをまとめて管理するための仕組み
- **ReplicaSet**: 同じ内容のPodを指定個数立てる
- **Deployment**: ReplicaSetを使い安全にバージョン移行



# ReplicaSet / Deployment

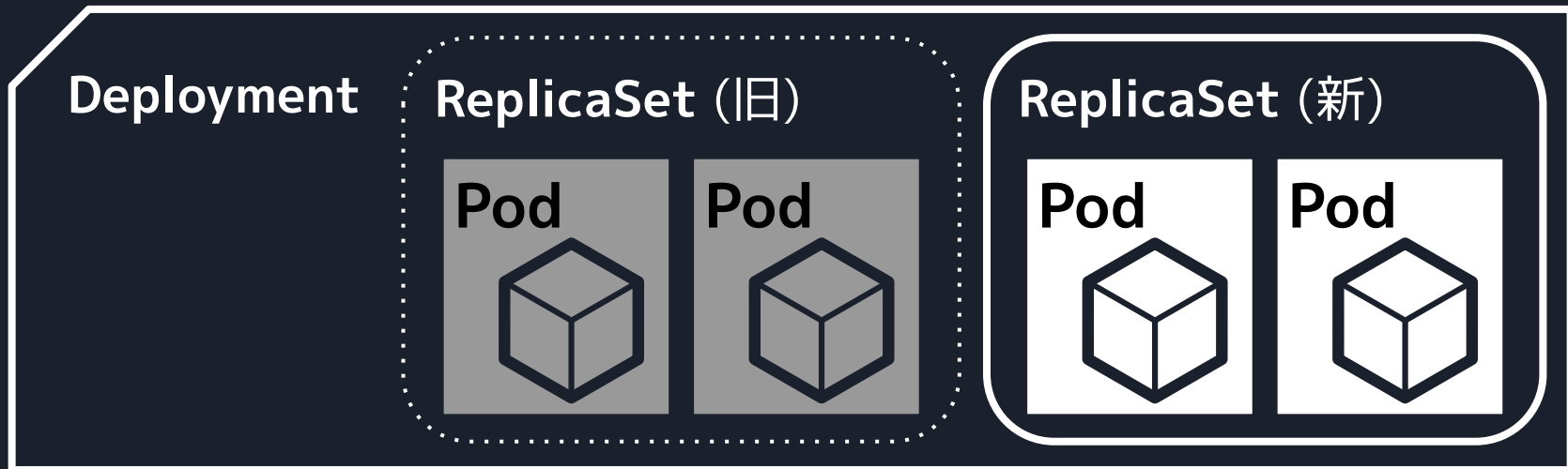
- Podをまとめて管理するための仕組み
- **ReplicaSet**: 同じ内容のPodを指定個数立てる
- **Deployment**: ReplicaSetを使い安全にバージョン移行





# ReplicaSet / Deployment

- Podをまとめて管理するための仕組み
- **ReplicaSet**: 同じ内容のPodを指定個数立てる
- **Deployment**: ReplicaSetを使い安全にバージョン移行



# ReplicaSet / Deployment

- Podをまとめて管理するための仕組み
- **ReplicaSet**: 同じ内容のPodを指定個数立てる
- **Deployment**: ReplicaSetを使い安全にバージョン移行

Deployment

ReplicaSet (新)

Pod



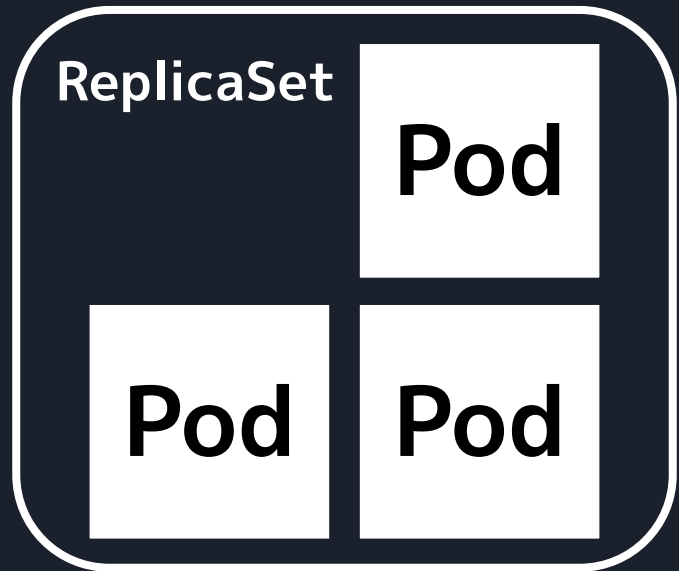
Pod





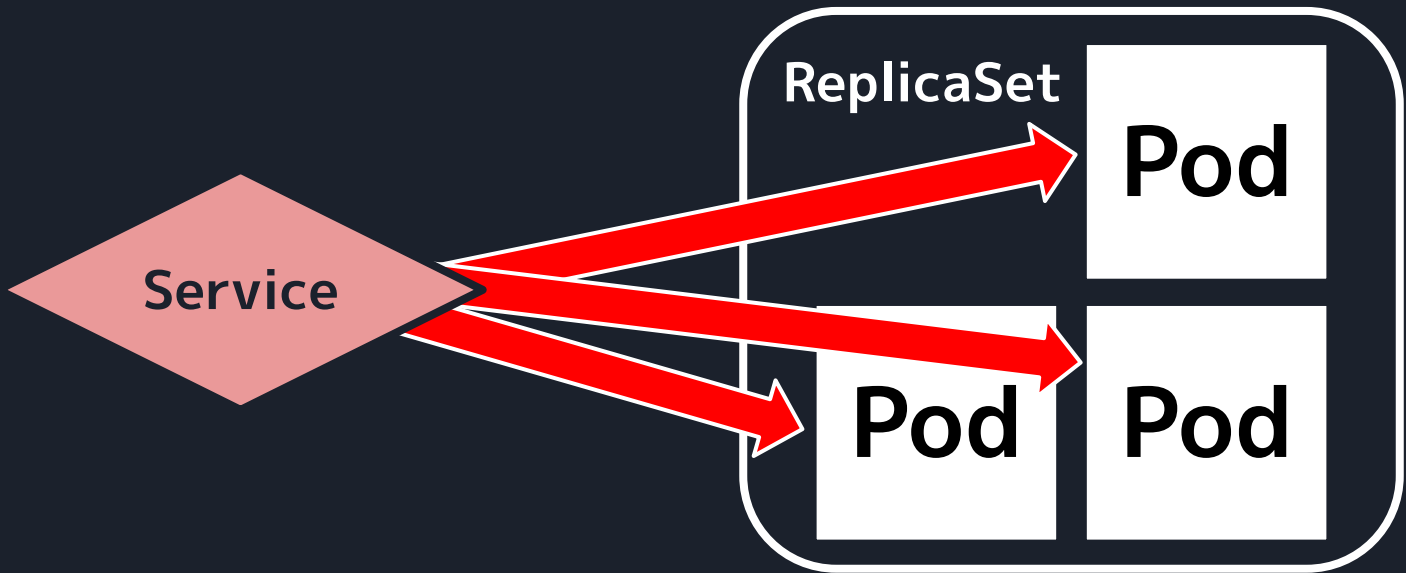
# Service

- 同じ内容のPod群に、均等に**通信を振り分ける** (L4LB)
- Pod群を**1つのIPアドレスでまとめる**ことができる



# Service


- 同じ内容のPod群に、均等に**通信を振り分ける** (L4LB)
- Pod群を**1つのIPアドレス**でまとめることができる





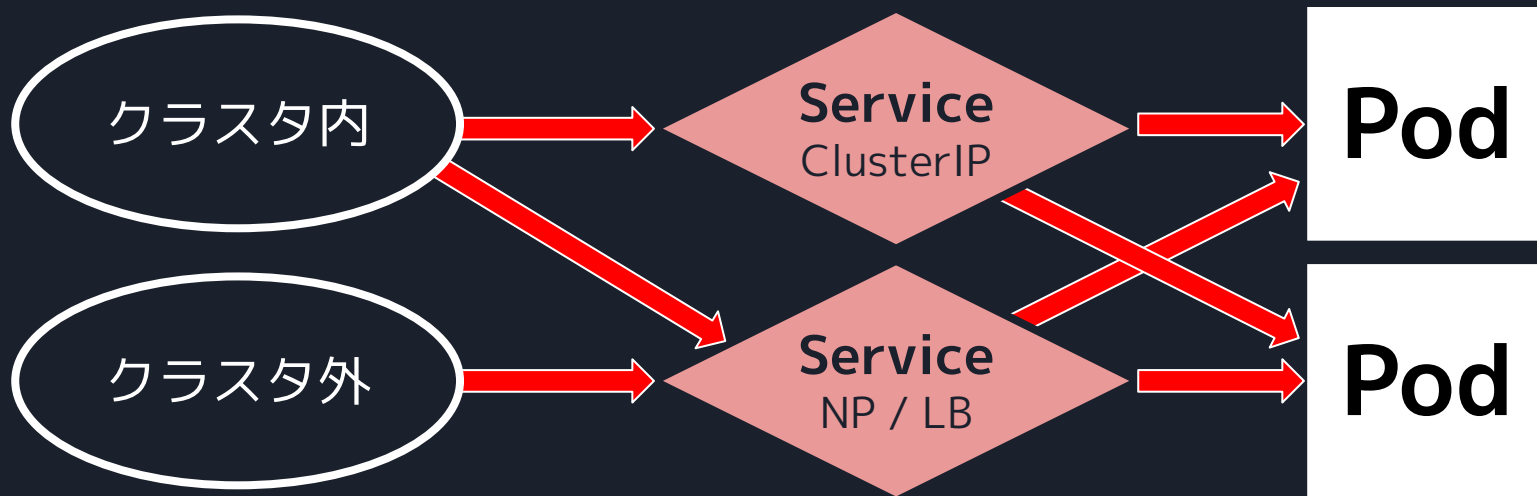
## Service

Podの入れ替わりで  
各Podの**IPが変わっても**  
**安定した仮想IPで**  
アクセスできる

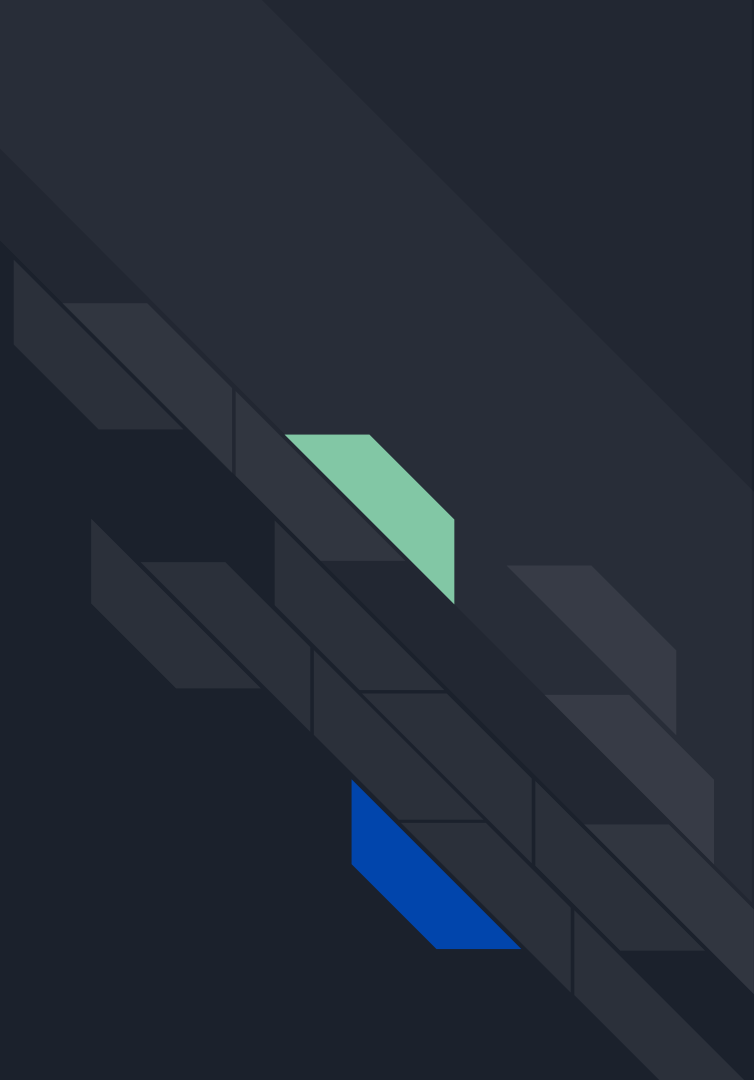


# Serviceのタイプ

- ClusterIP: **クラスタ内**からの通信のみ受け付ける
- NodePort / LoadBalancer: **クラスタ外**からの通信も受け付ける

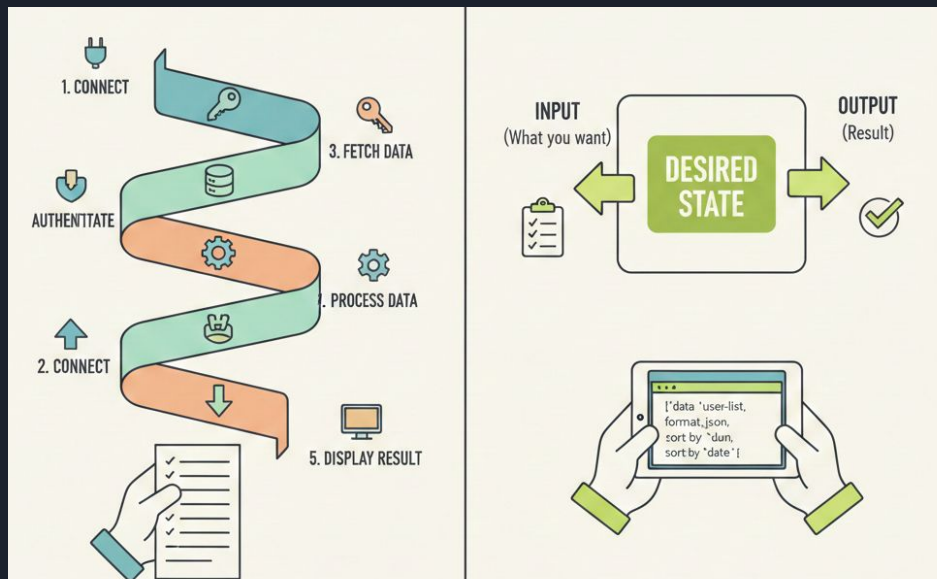


Kubernetesの実態:  
**宣言的API と**  
**Reconciliation Loop**



# 宣言的なAPI

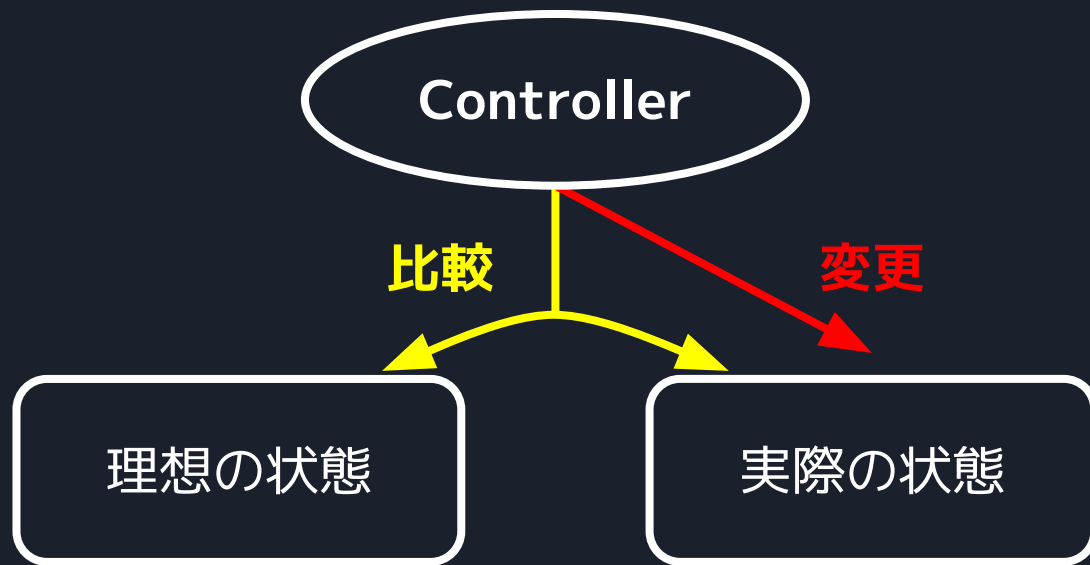
- 「これをして下さい」でなく「この状態にして下さい」
  - 「マニフェスト」で理想の状態を書く
- 暗黙知のない  
安定した稼働





# Reconciliation Loop

- リソースを理想の状態に持っていくための操作



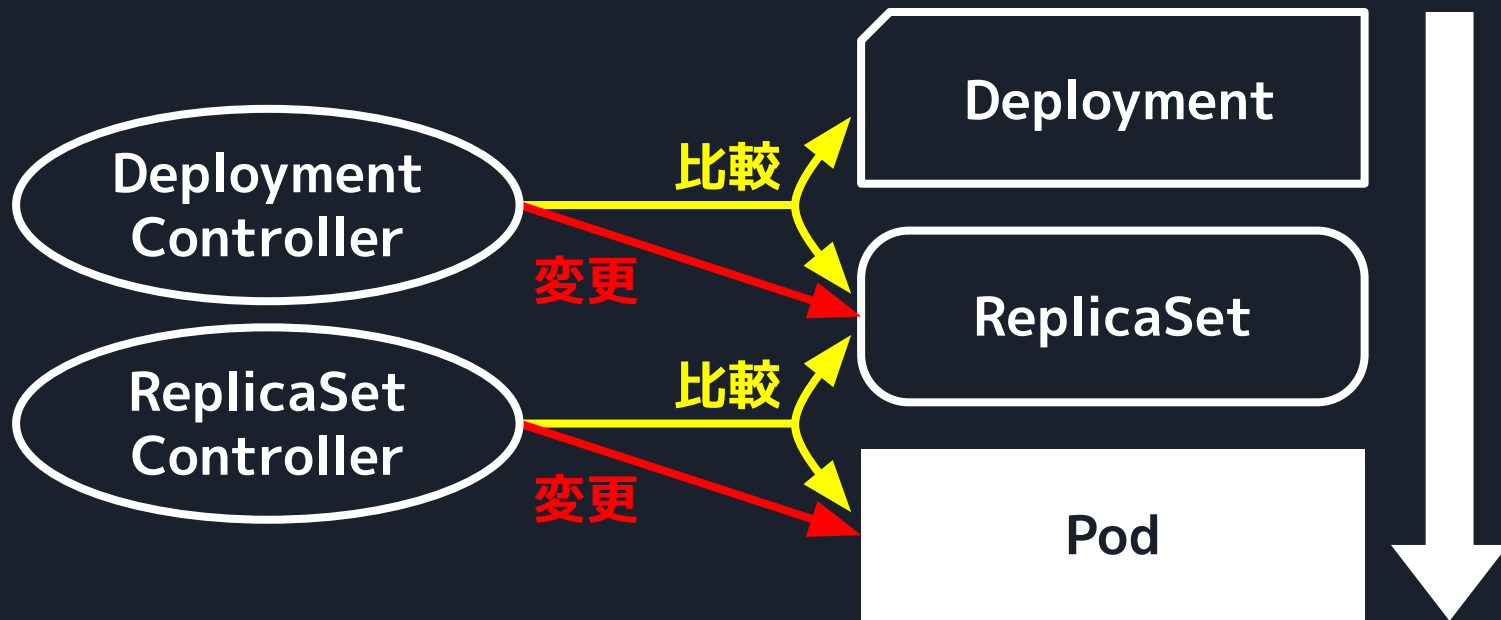



# Reconciliation Loop

- 以下の操作を一定時間ごとに繰り返す
  - **理想の状態** (宣言されたマニフェスト) を取得
  - **実際の状態**を観測
  - 理想の状態と実際の状態を**比較**
  - 2つの差を埋めるように、**実際の状態を変更**する
- **Controller**と呼ばれる部品がこれを担当する
- **Self-Healing**が簡単に実装できる

# 例: DeploymentがPodになるまで

- 単純なReconciliation Loopの積み重ねで複雑性を実現





# 宣言された理想の状態を保存し Reconciliation Loopを回す ただそれだけのシステム

ベースのアイデアは非常に単純



# 覚えておいてもらいたいこと

- K8sは拡張性が高いが、**ネットワークの視点**で見たとき重要なリソースは**少ない**
  - ほとんどの内部通信は、**PodからServiceを經由しPodへ**流れる
- **宣言的APIとReconciliation Loop**がK8sの核
  - K8sのネットワークもこの二つに準じて設定される必要がある
  - **手続き的な物を宣言的にする実装が必要**



# アジェンダ

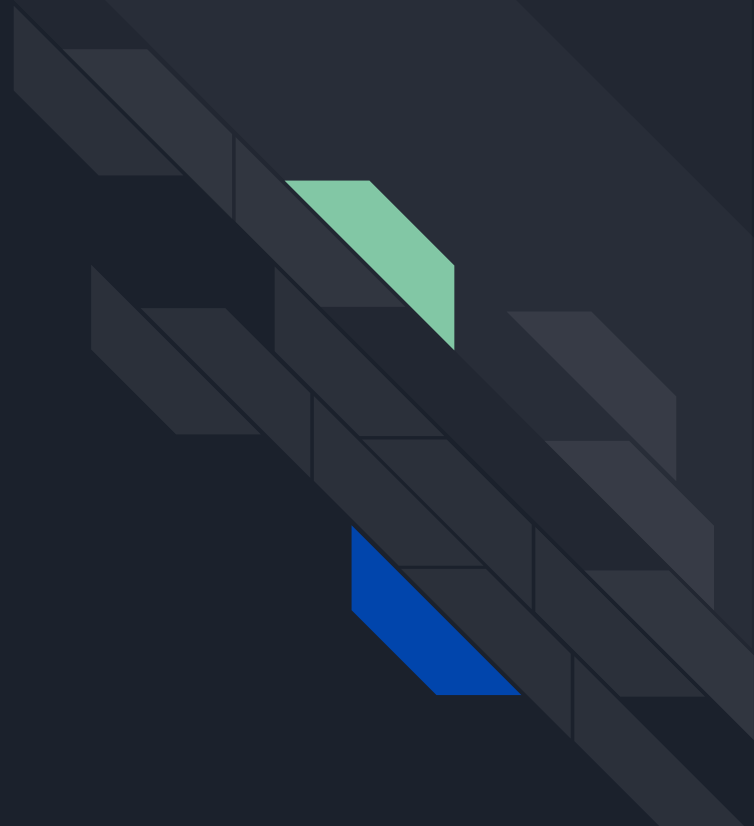
- イントロダクション
- Kubernetesを知る ←イマココ
- Kubernetesとネットワーク
- Linuxのネットワークスタック
- L2/L3ネットワーク: CNI
- L4ロードバランサー: kube-proxy
- プラットフォームとの対話で広がる可能性
- まとめ



# アジェンダ

- イントロダクション
- Kubernetesを知る
- Kubernetesとネットワーク ←イマココ
- Linuxのネットワークスタック
- L2/L3ネットワーク: CNI
- L4ロードバランサー: kube-proxy
- プラットフォームとの対話で広がる可能性
- まとめ

# Kubernetesと ネットワーク





# Kubernetesを流れる通信 3種

- Pod → Service → Pod
  - クラスタ内通信
  - ノード内とノード間を考慮する必要がある



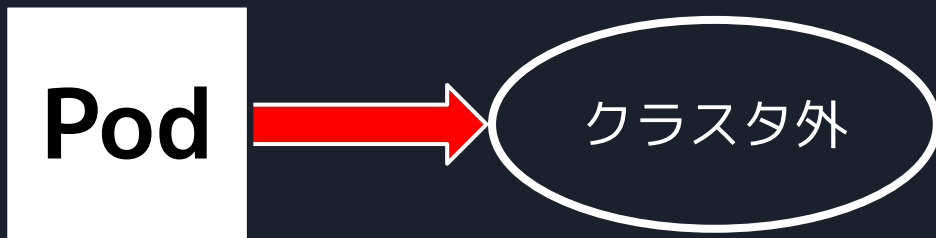
# Kubernetesを流れる通信 3種

- 外部 → Service → Pod
  - クラスタ外からのIngress通信
  - **Gateway API** (L7LB) というリソースとの関わりが深いが、本プログラムでは割愛



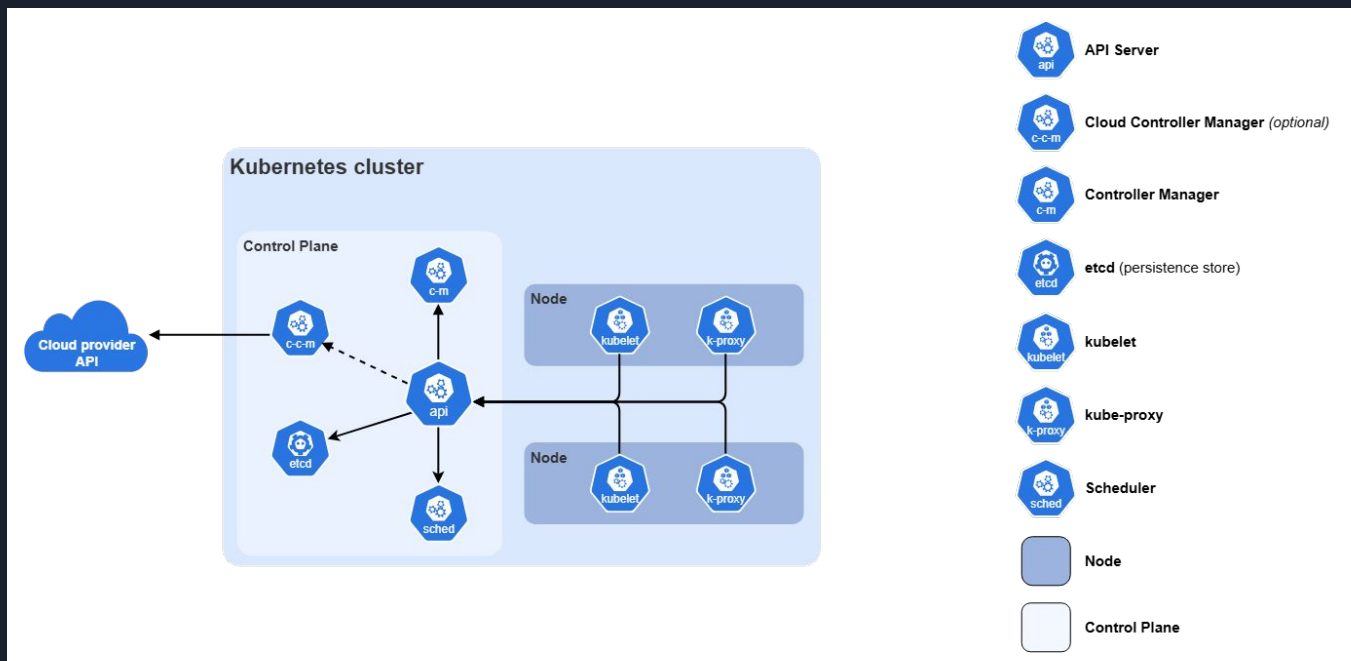
# Kubernetesを流れる通信 3種

- Pod → 外部
  - クラスタ外へのEgress通信
  - Podからのインターネット疎通を保証する



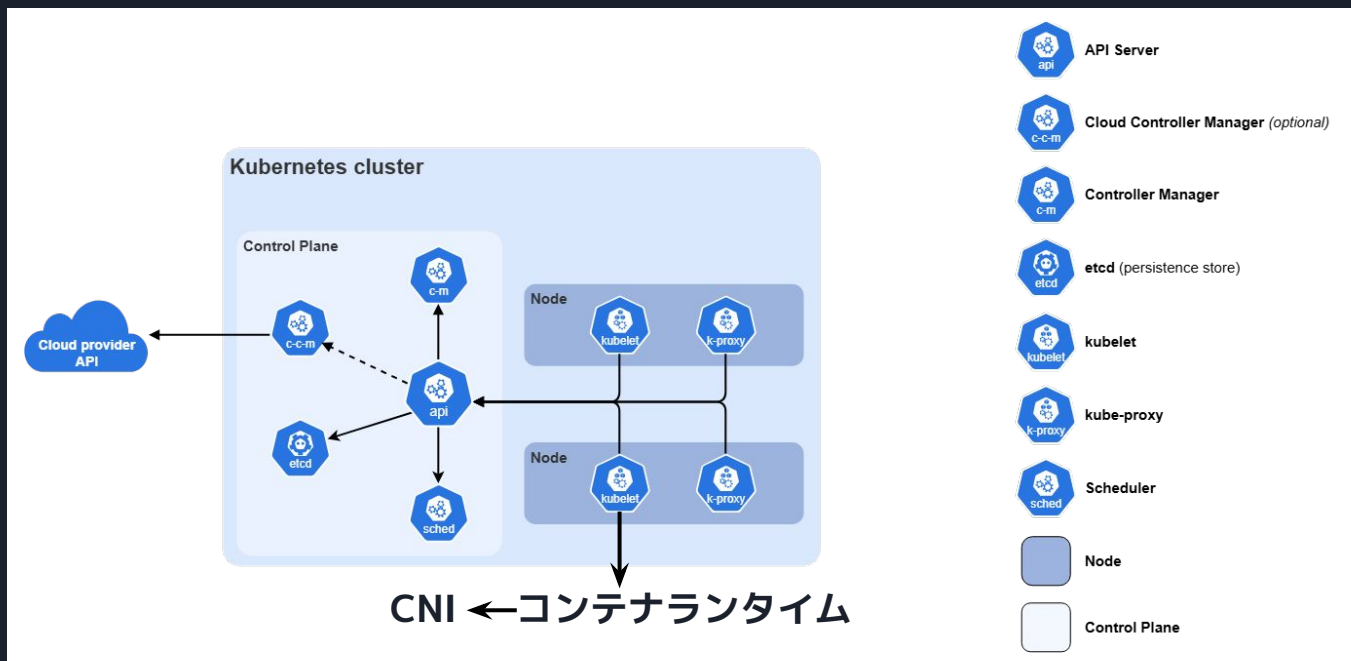
# Kubernetesを構成する部品たち

<https://kubernetes.io/docs/concepts/overview/components/> より引用



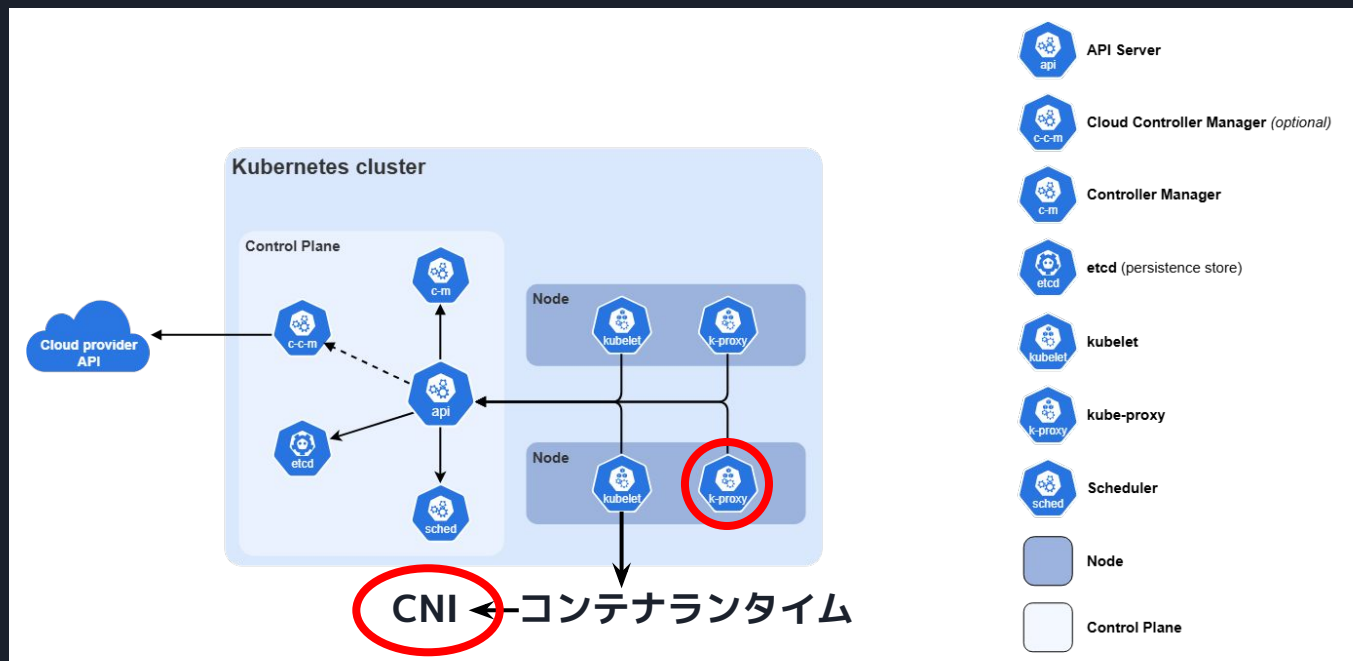
# Kubernetesを構成する部品たち

- K8sがデフォルトで用意していないものもある



# ネットワークを担当する部品たち

- CNI と kube-proxy



# ネットワークを担当する部品たち

- CNI

- Pod → Pod / Pod → 外部 を担当
- PodのL2/L3ネットワーク疎通を担保する
- Kubernetesのデフォルト実装はない



- kube-proxy

- Service → Pod / 外部 → Service を担当
- 複数Podに対するL4LB (Service) を実装する
- Kubernetesがデフォルト実装を用意している





ここからは  
それぞれの部品の**詳細と**  
**実装**を見ていきます

どのように「架け橋」が用意されているかにご注目下さい





# アジェンダ

- イントロダクション
- Kubernetesを知る
- Kubernetesとネットワーク ←イマココ
- Linuxのネットワークスタック
- L2/L3ネットワーク: CNI
- L4ロードバランサー: kube-proxy
- プラットフォームとの対話で広がる可能性
- まとめ




# アジェンダ

- イントロダクション
- Kubernetesを知る
- Kubernetesとネットワーク
- Linuxのネットワークスタック ←イマココ
- L2/L3ネットワーク: CNI
- L4ロードバランサー: kube-proxy
- プラットフォームとの対話で広がる可能性
- まとめ

# Linuxの ネットワークスタック





それぞれの詳細を見る前に  
まずは**前提知識**となる  
**Linux**のネットワークをおさらい

LinuxのSDNを一度俯瞰してみましょう



# Network Namespace (netns)

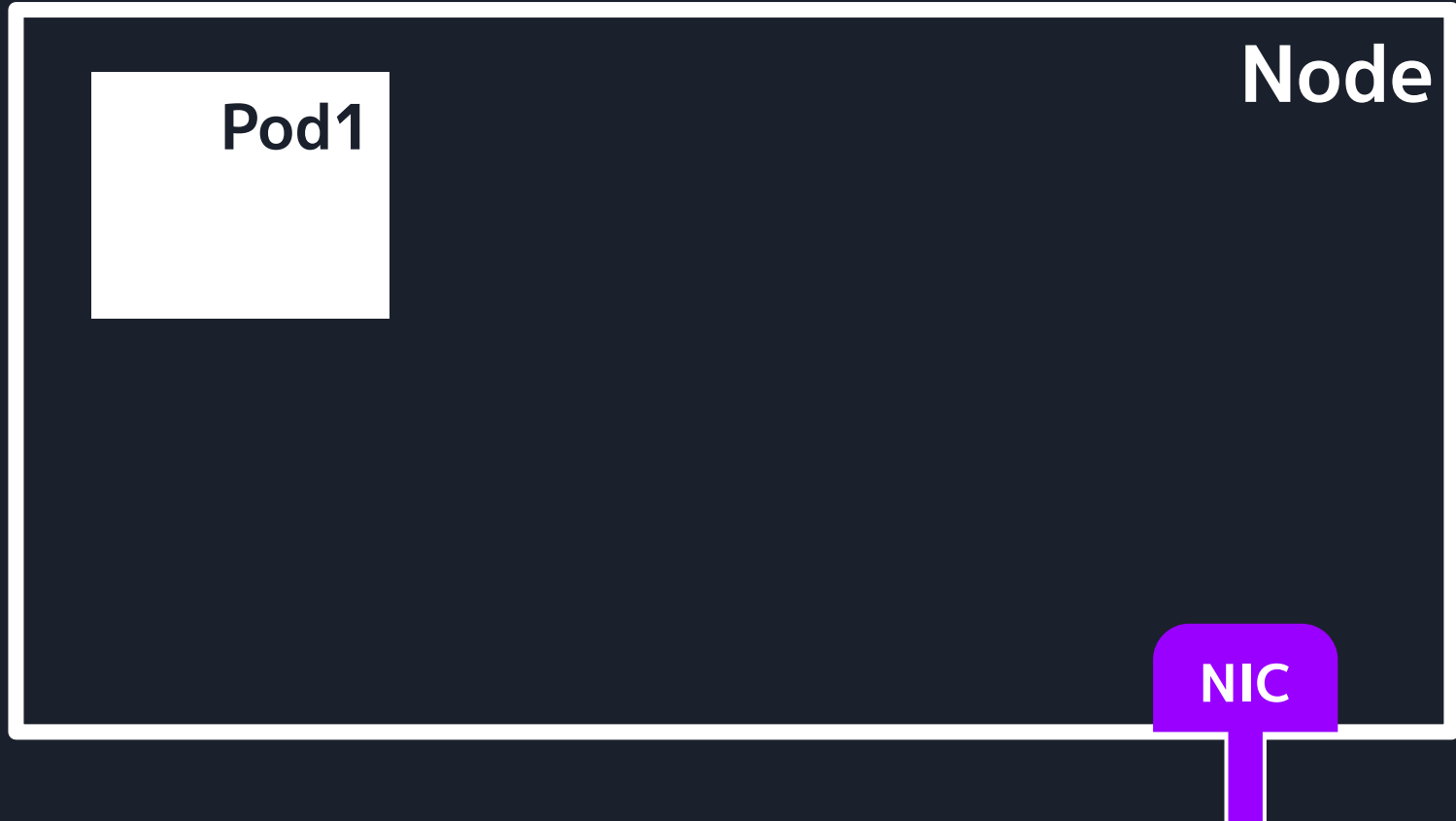
- Linuxの機能の一つ
- 1台のマシンの中に**独立した仮想のホスト**を作る
- **netns1つにつき、1マシン**と考えてよい
  - **Podの** (ネットワーク的な) **実体**
  - それぞれ個別に、NICを生やしたりIP/MACアドレスを持たせたりできる



# Network Namespace (netns)



# Network Namespace (netns)

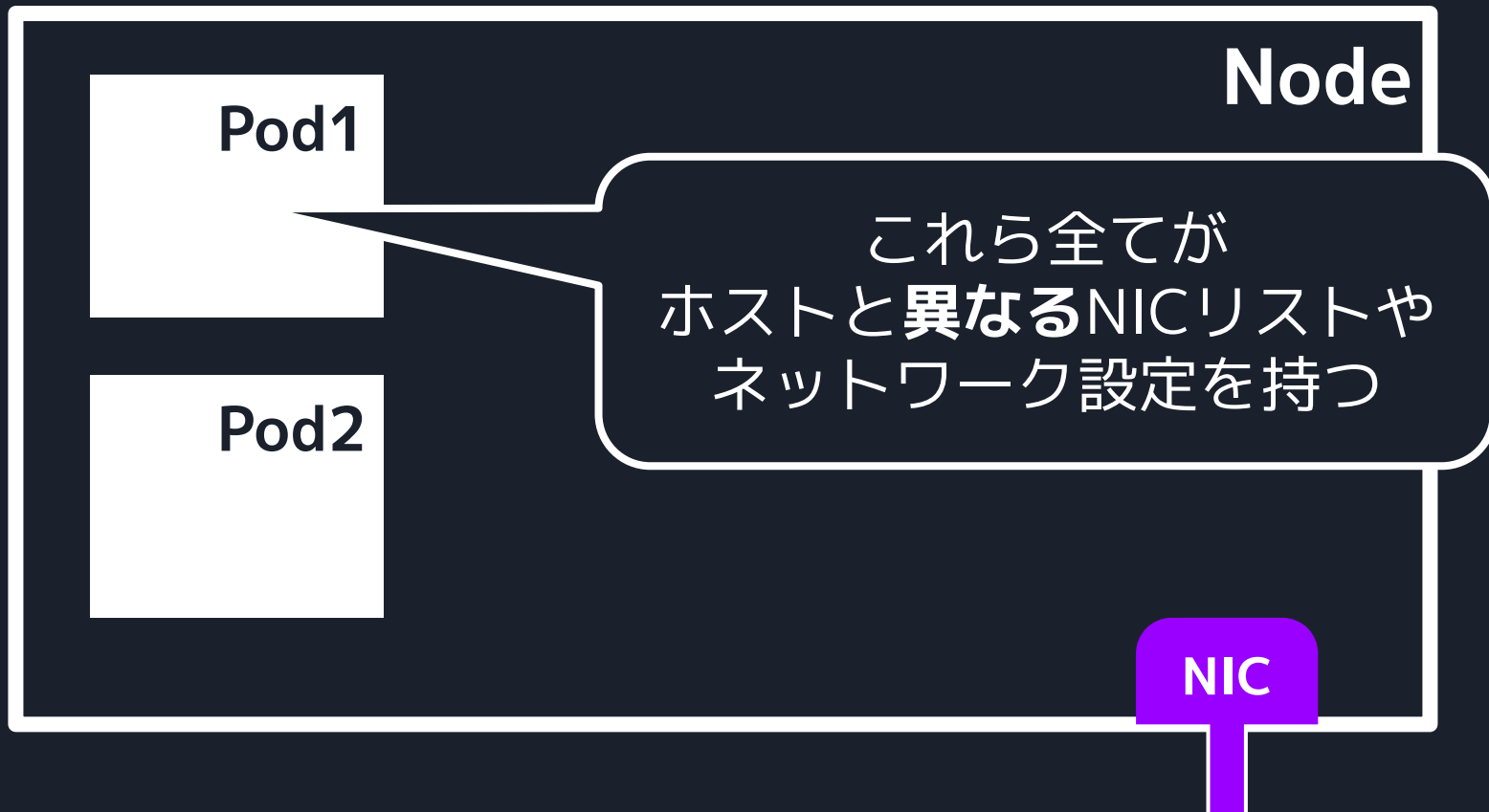


# Network Namespace (netns)





# Network Namespace (netns)



実は、親マシンのネットワークも…



実は、親マシンのネットワークも…

親マシンというHWの**枠**に  
(親マシン**自身**のNWも含め)  
**区切られたネットワーク**が  
たくさんあるイメージ

NIC

実は、親マシ

ハードウェアの境界線と  
ネットワークの境界線が  
別物になるということ

親マシン  
(親マシン

区切られたネットワークが  
たくさんあるイメージ

NIC

# veth

- Linux内で**仮想的**に作られるイーサネットケーブル
  - 繋がった**NICのペア**として生成される
- これによりnetns同士を接続することが可能になる



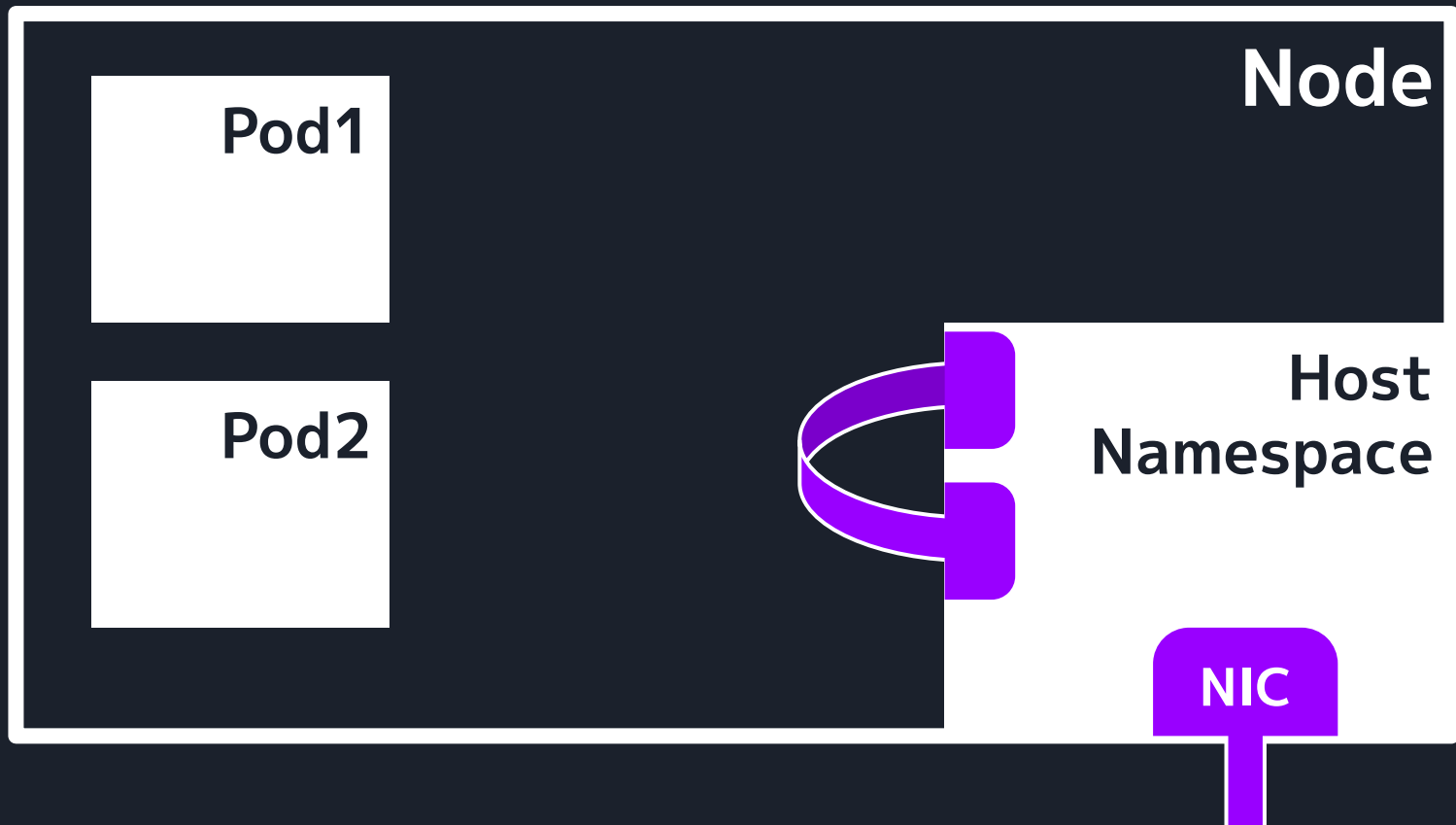


# veth

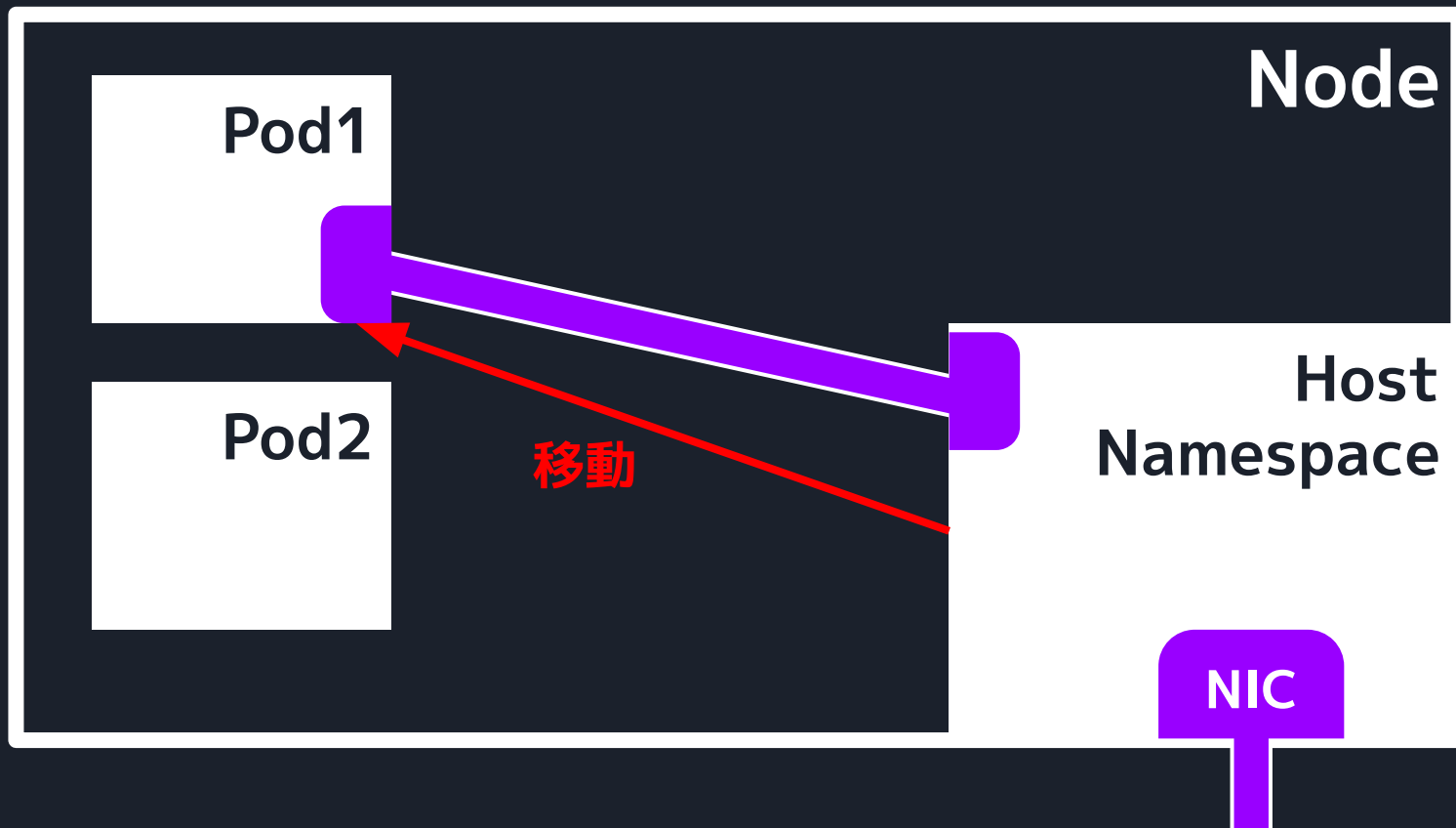




# veth



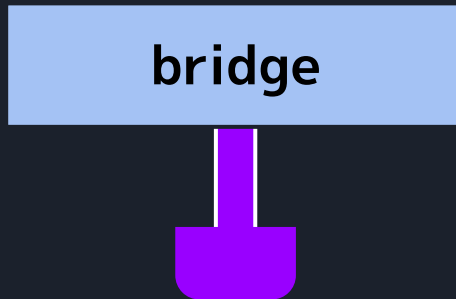
veth



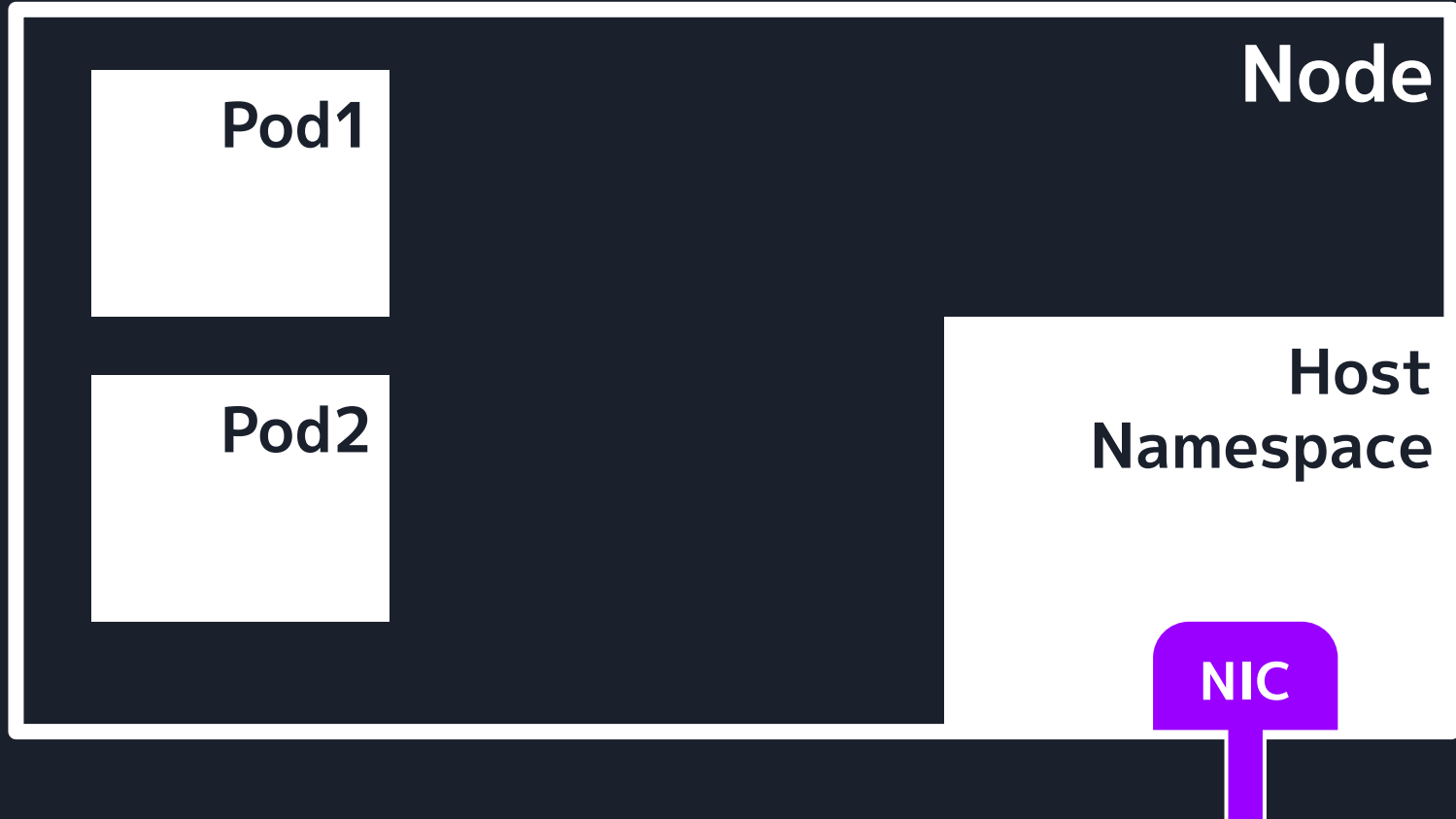


# Linux Bridge

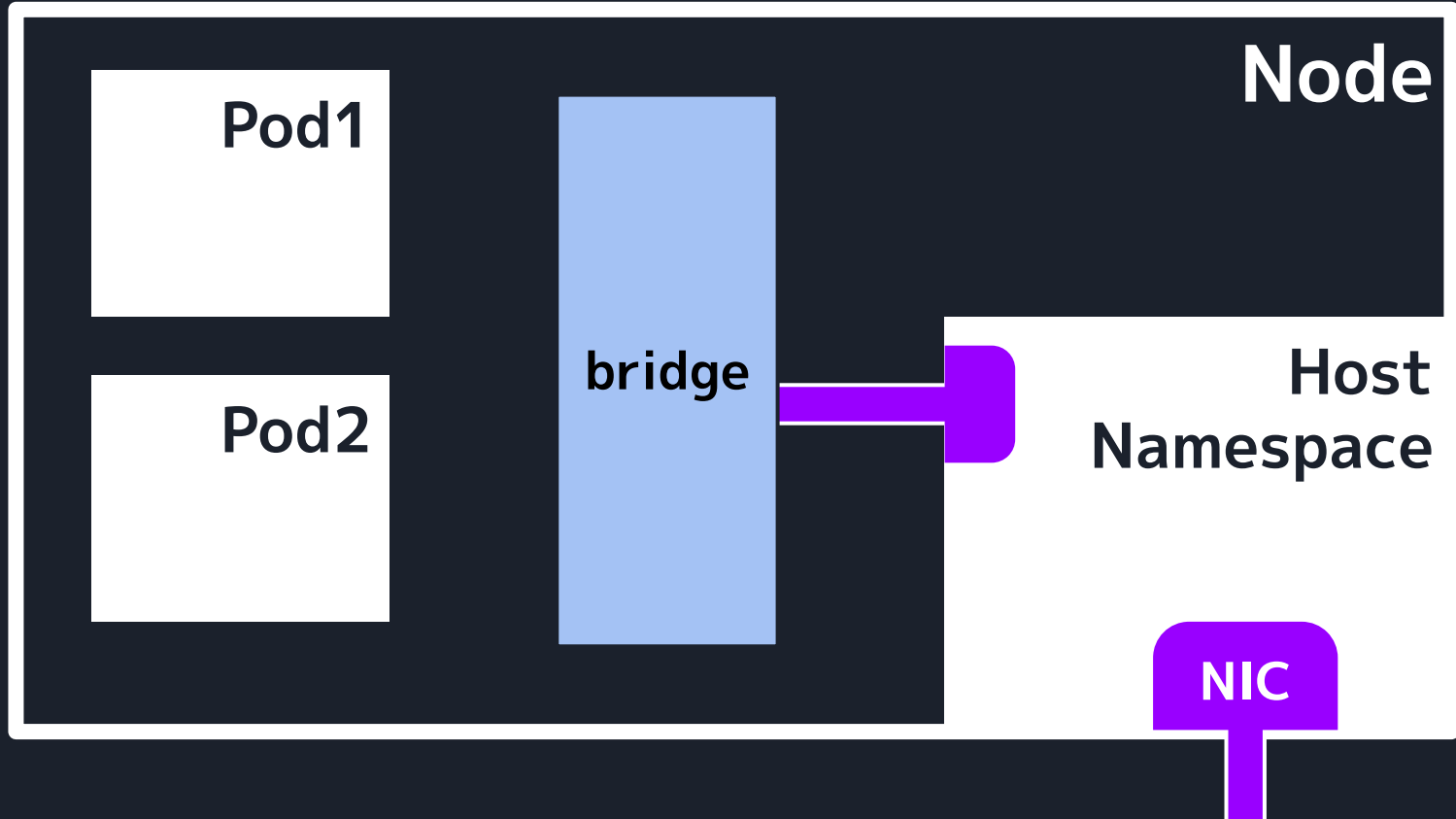
- 仮想L2スイッチ
  - 複数のvethを接続すると、接続されたnetns同士のEthernet疎通・ARP解決ができるようになる
- Linuxの都合で、host namespaceと接続されたNICがあらかじめ生成される



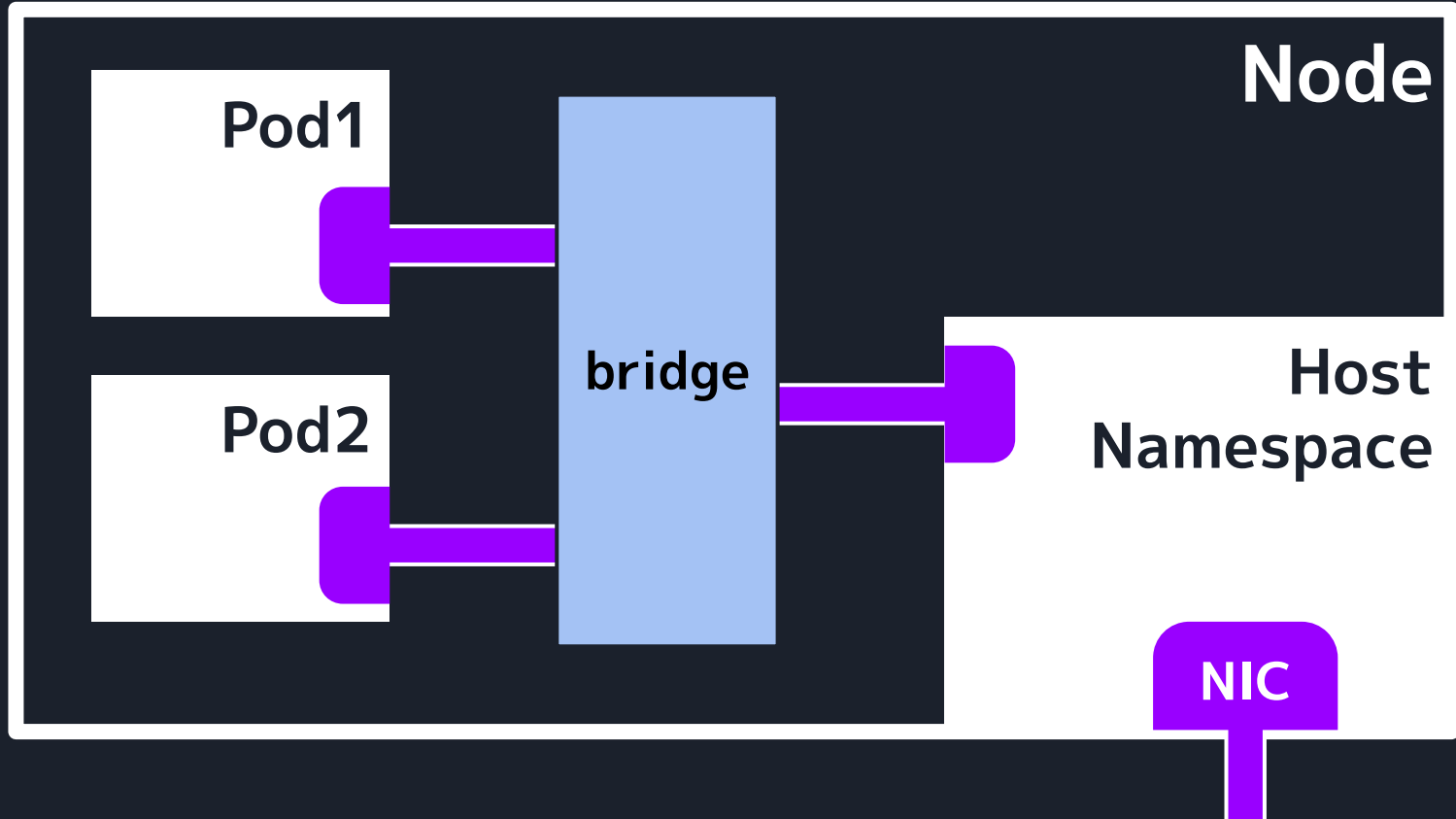
# Linux Bridge




# Linux Bridge



# Linux Bridge

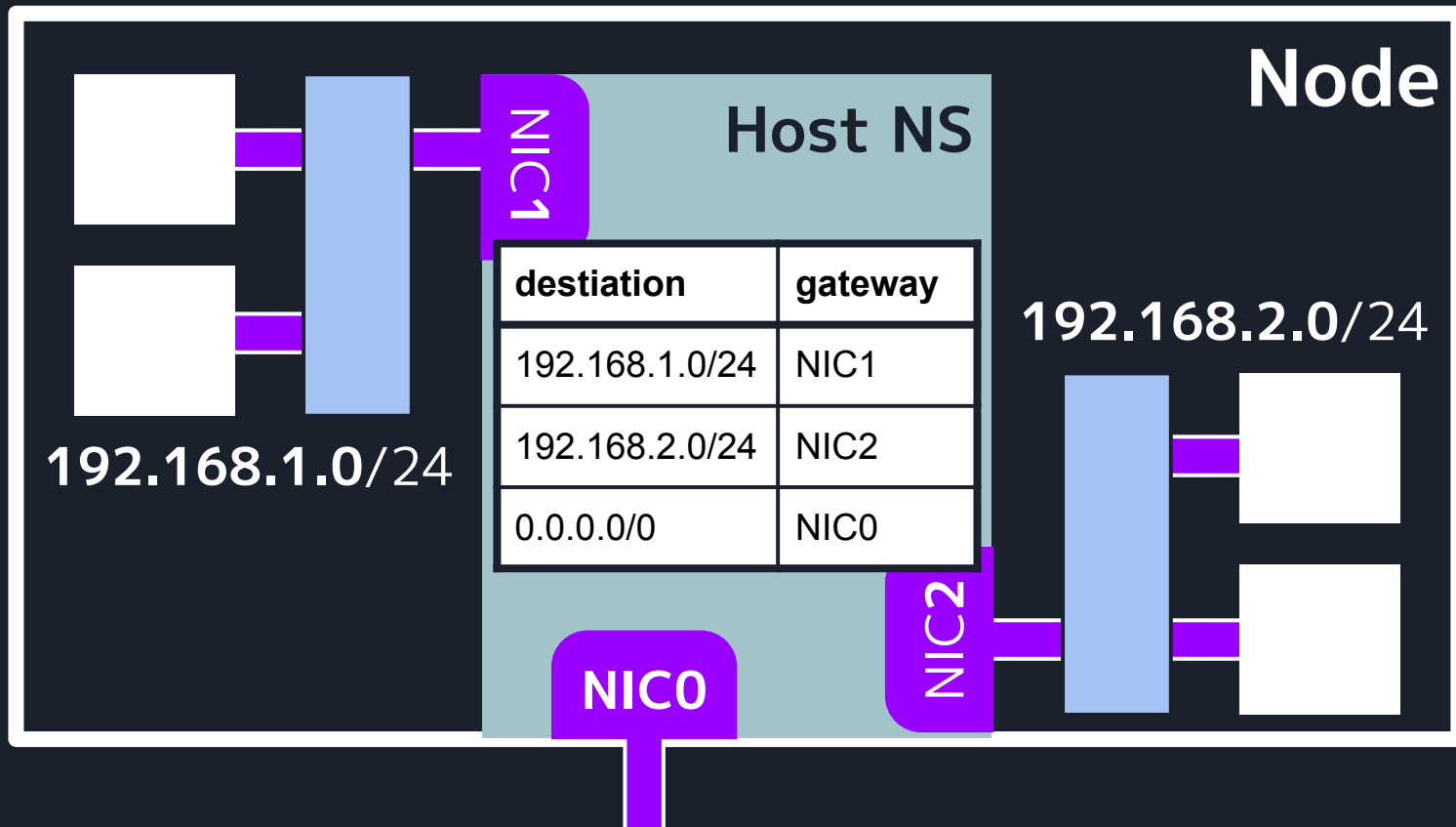




ここまでの機能を使うと  
1台の親マシンの中で**複雑な  
サブネット構造**も作れる

たくさんのルーターが複雑に絡んだ構成も再現可能

## ex) 複雑なサブネット構造





# netfilter

- パケットをフィルタ・加工・転送できる仕組み
  - Filter
    - 特定のIPアドレス・ポート宛のパケットを落とす
  - NAT
    - 送り元・送り先のIPアドレス・ポートを変更する
- 操作方法 (フロントエンド)
  - 長らくiptablesを使って設定されていたが、**非推奨**に
  - 後継のnftablesへの移行が進んでいる



# アジェンダ

- イントロダクション
- Kubernetesを知る
- Kubernetesとネットワーク
- Linuxのネットワークスタック ←イマココ
- L2/L3ネットワーク: CNI
- L4ロードバランサー: kube-proxy
- プラットフォームとの対話で広がる可能性
- まとめ

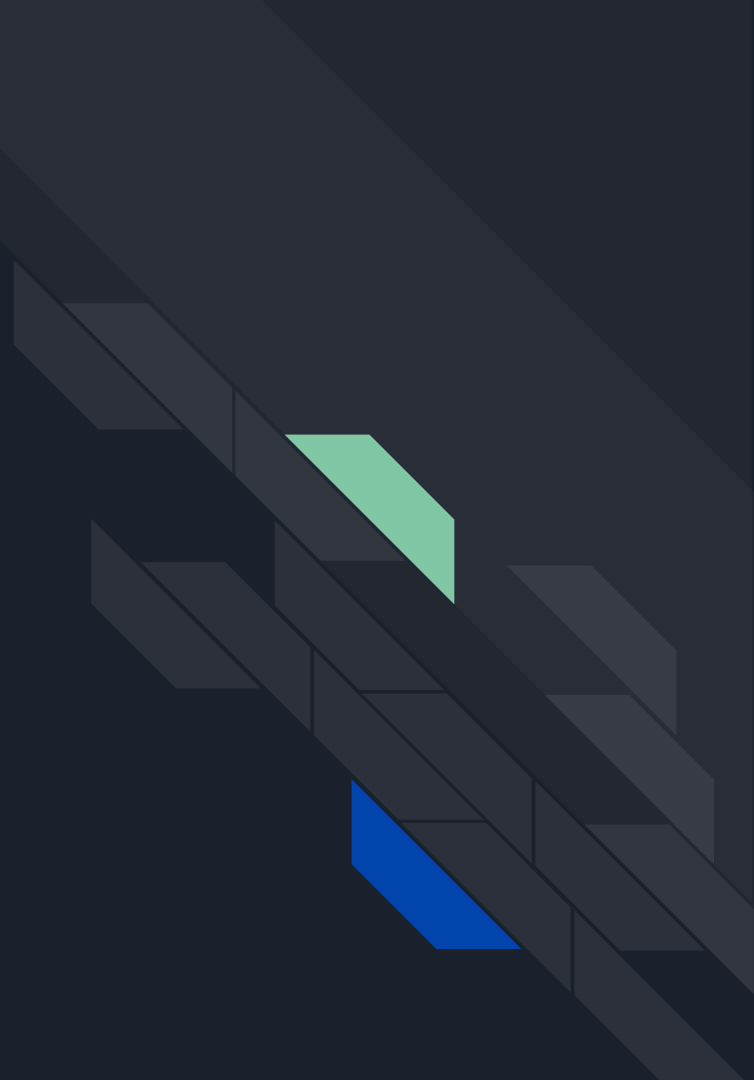




# アジェンダ

- イントロダクション
- Kubernetesを知る
- Kubernetesとネットワーク
- Linuxのネットワークスタック
- L2/L3ネットワーク: CNI ←イマココ
- L4ロードバランサー: kube-proxy
- プラットフォームとの対話で広がる可能性
- まとめ

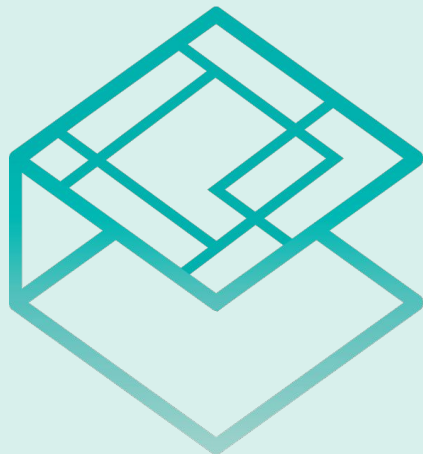
# L2/L3ネットワーク: CNI





# CNI

- Container Network Interface
- Podを何らかのネットワークに参加させるための仕組み



CNI












# CNIはKubernetesだけじゃない！

- 意外と多くのコンテナ基盤がCNIにネットワーク機能を移譲している

## Who is using CNI?

### Container runtimes

- [Kubernetes](#)  - a system to simplify container operations
- [HashiCorp Nomad](#)  - A simple and flexible scheduler and orchestrator to deploy and manage containers and non-containerized applications across on-prem and clouds at scale.
- [Containerd](#)  - A CRI-compliant container runtime
- [cri-o](#)  - A lightweight container runtime
- [OpenShift](#)  - Kubernetes with additional enterprise features
- [Cloud Foundry](#)  - a platform for cloud applications
- [Apache Mesos](#)  - a distributed systems kernel
- [Amazon ECS](#)  - a highly scalable, high performance container management service
- [Singularity](#)  - a container platform optimized for HPC, EPC, and AI
- [OpenSVC](#)  - an orchestrator for legacy and containerized application stacks

# CNIはKubernetesだけじゃない！

- 意外と多くのコンテナ  
移譲している

Containerd (ただの  
コンテナランタイム)  
からもCNIは使える！

Who is using CNI?

## Container runtimes

- [Kubernetes](#) - a system to simplify container operations
- [HashiCorp Nomad](#) - A simple and flexible scheduler and orchestrator to deploy and manage containers and non-containerized applications across on-prem and clouds at scale.
- [Containerd](#) - A CRI-compliant container runtime
- [cri-o](#) - A lightweight container runtime
- [OpenShift](#) - Kubernetes with additional enterprise features
- [Cloud Foundry](#) - a platform for cloud applications
- [Apache Mesos](#) - a distributed systems kernel
- [Amazon ECS](#) - a highly scalable, high performance container management service
- [Singularity](#) - a container platform optimized for HPC, EPC, and AI
- [OpenSVC](#) - an orchestrator for legacy and containerized application stacks



## 余談: CNIは (実質) Kubernetesだけ

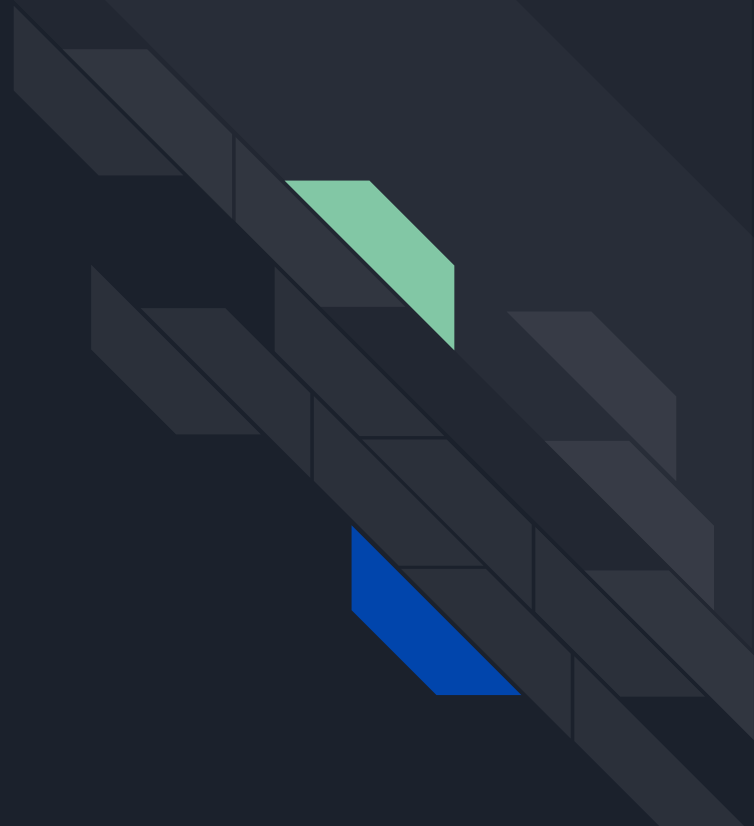
- DockerやPodmanなど、多くの人がローカルで使うコンテナランタイムは、独自の**ドライバ**を使って管理する
  - Containerd (Dockerの中身) を直接触る人は少ない
- 実際のユースケースは、ほぼ**Kubernetesが類似のオーケストレーションサービス**
- Kubernetesのネットワークのメンテナーも「CNIはK8sに**特化した仕組みじゃないのに**、ほぼK8sからしか**使われない**」と言っている



# 広義「CNI」とは

- **CNI仕様** (containernetworkinterface/cni)
  - コンテナ基盤がどのようにプラグインを呼び出すかの定義
- **基盤の要件定義**
  - 基盤としてどんなネットワークが欲しいのか
- **CNIプラグイン**
  - CNI仕様を通して操作が可能で、基盤の要件定義を満たす**実際のネットワーク実装**

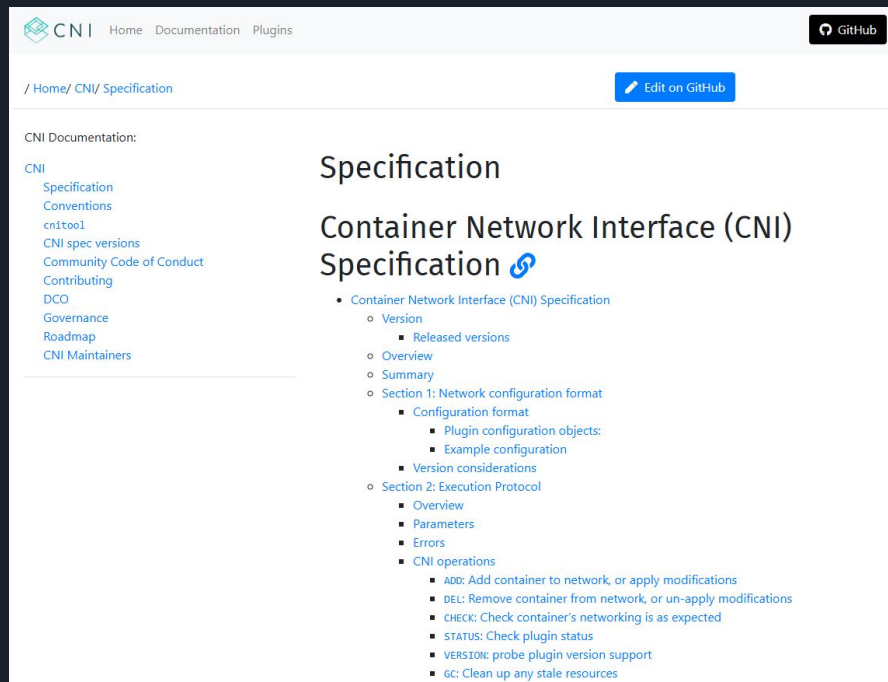
CNI仕様





# CNI仕様

- プラグインが満たすべきユーザーとのインターフェースが定められている
  - 設定方法
  - 操作方法
  - 複数CNIの連携



The screenshot shows the CNI Specification page. The header includes the CNI logo, navigation links (Home, Documentation, Plugins), and a GitHub link. The main content area is titled 'Specification' and 'Container Network Interface (CNI) Specification'. A table of contents on the left lists various documents. The right side shows a detailed table of contents for the 'Container Network Interface (CNI) Specification' document, including sections like 'Version', 'Overview', 'Summary', 'Section 1: Network configuration format', and 'Section 2: Execution Protocol'.

CNI Home Documentation Plugins GitHub

/ Home/ CNI/ Specification Edit on GitHub

CNI Documentation:

- CNI
  - Specification
  - Conventions
  - cnitool
  - CNI spec versions
  - Community Code of Conduct
  - Contributing
  - DCO
  - Governance
  - Roadmap
  - CNI Maintainers

## Specification

### Container Network Interface (CNI) Specification

- Container Network Interface (CNI) Specification
  - Version
    - Released versions
  - Overview
  - Summary
  - Section 1: Network configuration format
    - Configuration format
      - Plugin configuration objects
      - Example configuration
    - Version considerations
  - Section 2: Execution Protocol
    - Overview
    - Parameters
    - Errors
    - CNI operations
      - ADD: Add container to network, or apply modifications
      - DEL: Remove container from network, or un-apply modifications
      - CHECK: Check container's networking is as expected
      - STATUS: Check plugin status
      - VERSION: probe plugin version support
      - GC: Clean up any stale resources



# CNI仕様 - 操作コマンド

- CNIプラグインは**5つのサブコマンド**を備えた**バイナリ**として実装される必要がある
  - ADDコマンド
  - DELコマンド
  - CHECKコマンド - 現在の状態を確認
  - GCコマンド - いらないリソースを掃除
  - VERSIONコマンド
- 特に大事なのが**ADD**と**DEL**

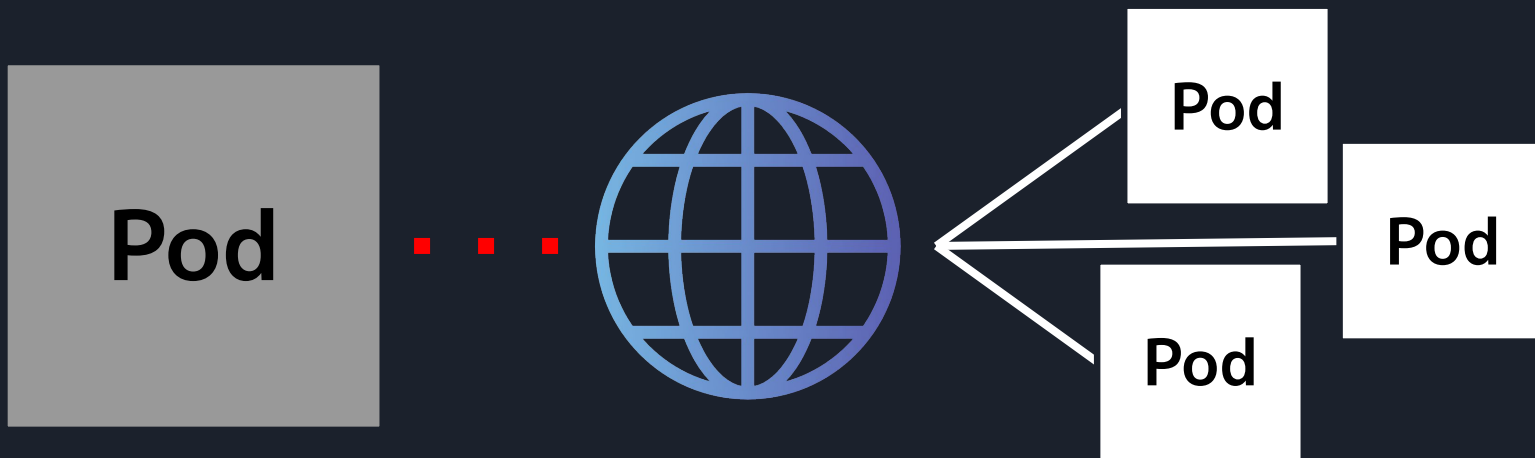
# CNI仕様 - 操作コマンド

- **ADD**コマンド
  - コンテナ (Pod) に**IPアドレス**を割り当てる
  - 指定されたコンテナを、用意された**コンテナネットワーク**の**ワーク**の中に**参加**させる



# CNI仕様 - 操作コマンド

- DELコマンド
  - 指定されたコンテナ (Pod) をネットワークから削除する

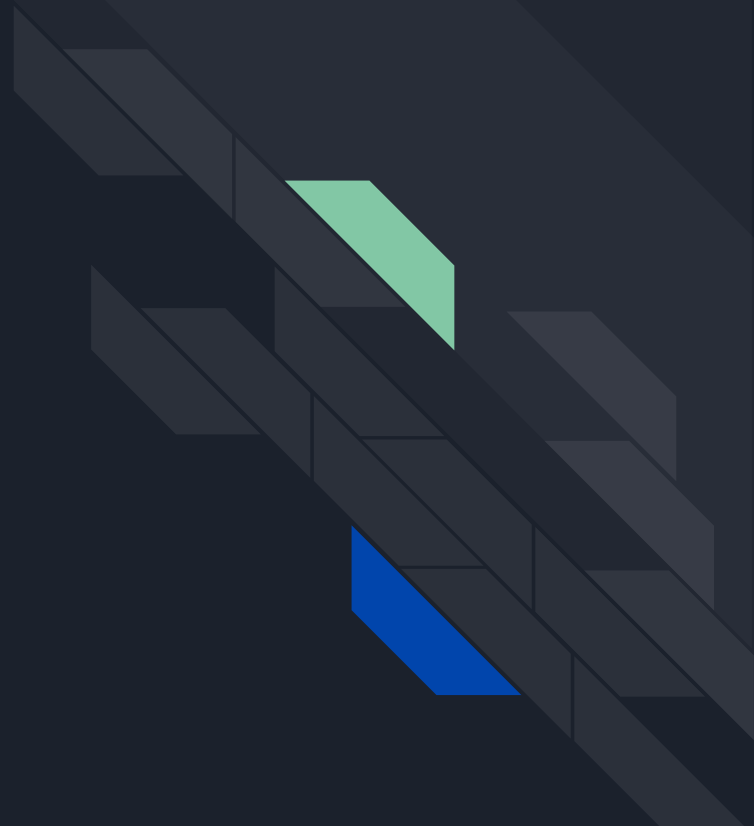




# CNI仕様

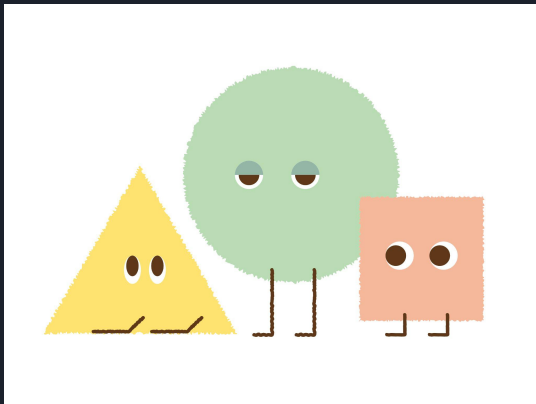
- 逆に言えば、この2つのコマンドが処理できればCNIプラグインとして動く
  - **Bash**でCNIプラグインを**作った**人もいます
- Kubernetesの場合CNIのコア機能を**デーモンで作動**させ**エントリーポイント**になるバイナリを別途用意し、そこからデーモンのAPIを叩きにいく構成が多い
  - 複雑なネットワーク**ライフサイクル**に対処するため

# 基盤の要件定義



# 基盤側の要件定義

- コンテナ基盤によって**必要な要件は違う**
- 例: **ローカル環境**
  - 同一ホスト内でコンテナ同士の疎通が取れる
  - コンテナから外部ネットワークに出られる





# The Kubernetes network model

- 各Podは**クラスタ内で唯一**のIPアドレスを持つ
- ホストの同じ/異なるに**関わらず**、**NAT無し**でPod同士の通信ができる

<https://kubernetes.io/docs/concepts/services-networking/>

## Services, Load Balancing, and Networking

Concepts and resources behind networking in Kubernetes.

The Kubernetes network model





# The Kubernetes network model

すなわち、CNIが**やらなければいけない**ことは

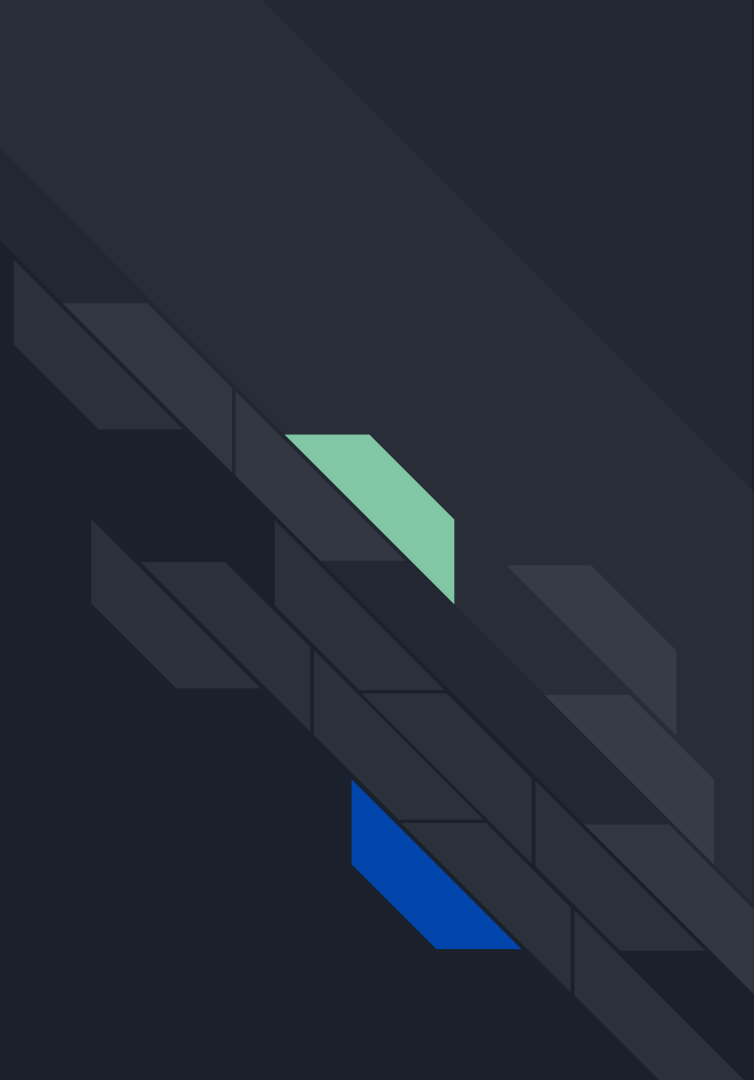
- Podにユニークな**IPアドレスを割り当てる** (IPAM)
- 同じクラスターのPod同士が、ホストの同じ/異なるに**関わらず、NATなし**で通信できるようにする



# CNIという仕様が担保する「自由」

- 逆に、これらが実現できれば使う手法は**なんでも良い**
  - これからいくつか例を見ていくが、**千差万別**
- CNI仕様は**どんなネットワーク技術でも受け入れる**広い受け皿
  - クラウド事業者が**元々のネットワーク基盤**を活用しやすいなど、様々なメリットがある

# CNIプラグインとその実装





# メジャーなOSSプラグイン

- **Flannel**
  - ノードを跨ぐPod間通信に**VXLAN**を使用したCNI
- **Calico**
  - 各ノードの持つPod CIDRを**BGP**で広報するCNI
- **Cilium**
  - ノード内の通信に**eBPF**を活用したCNI

この3つを例に、それぞれの**デフォルトの通信方式**がどのように疎通性を実現しているか見ていきましょう



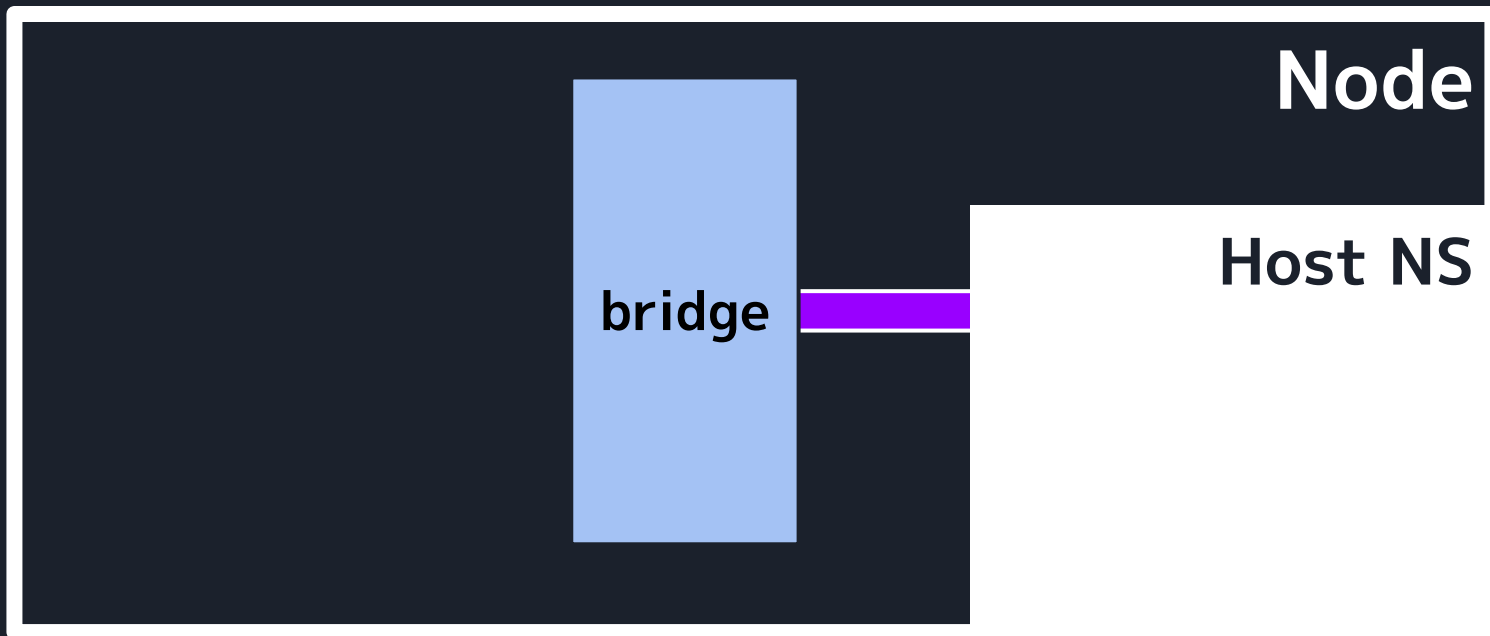
# Flannel - 同じノード内

下準備



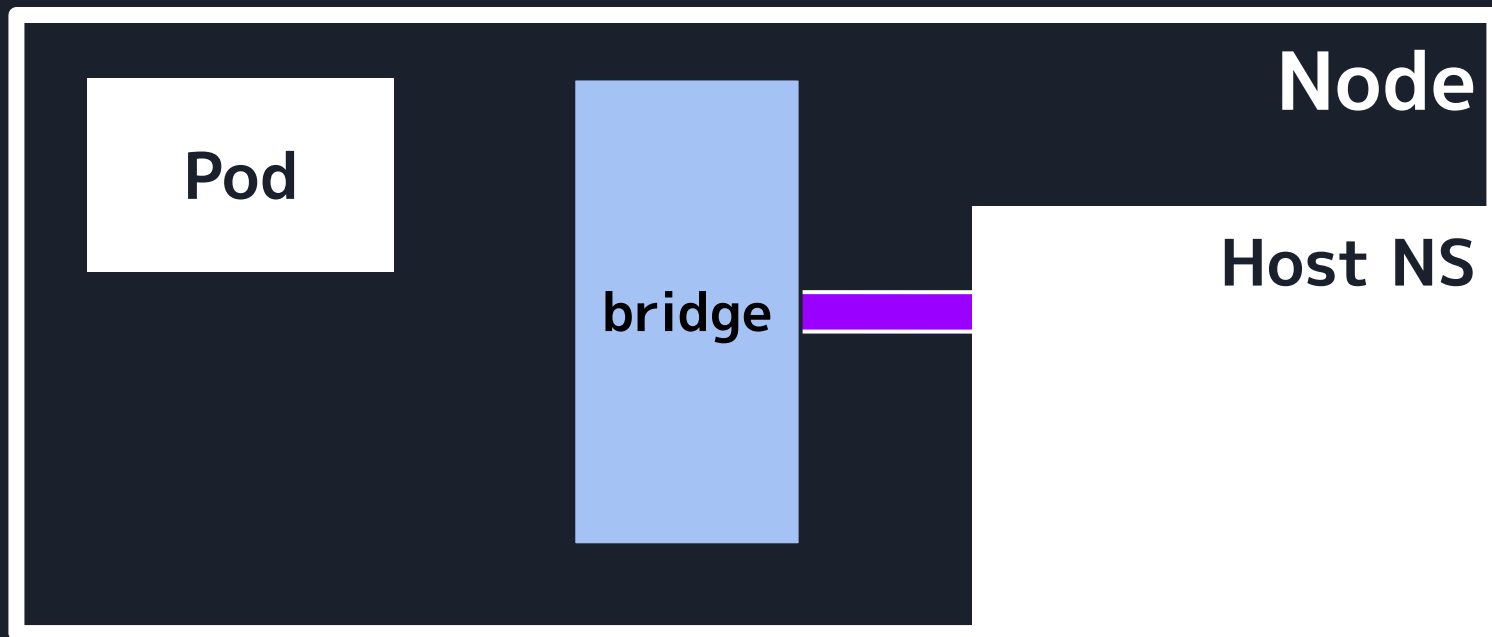
# Flannel - 同じノード内

下準備: **Linux Bridge**を作る



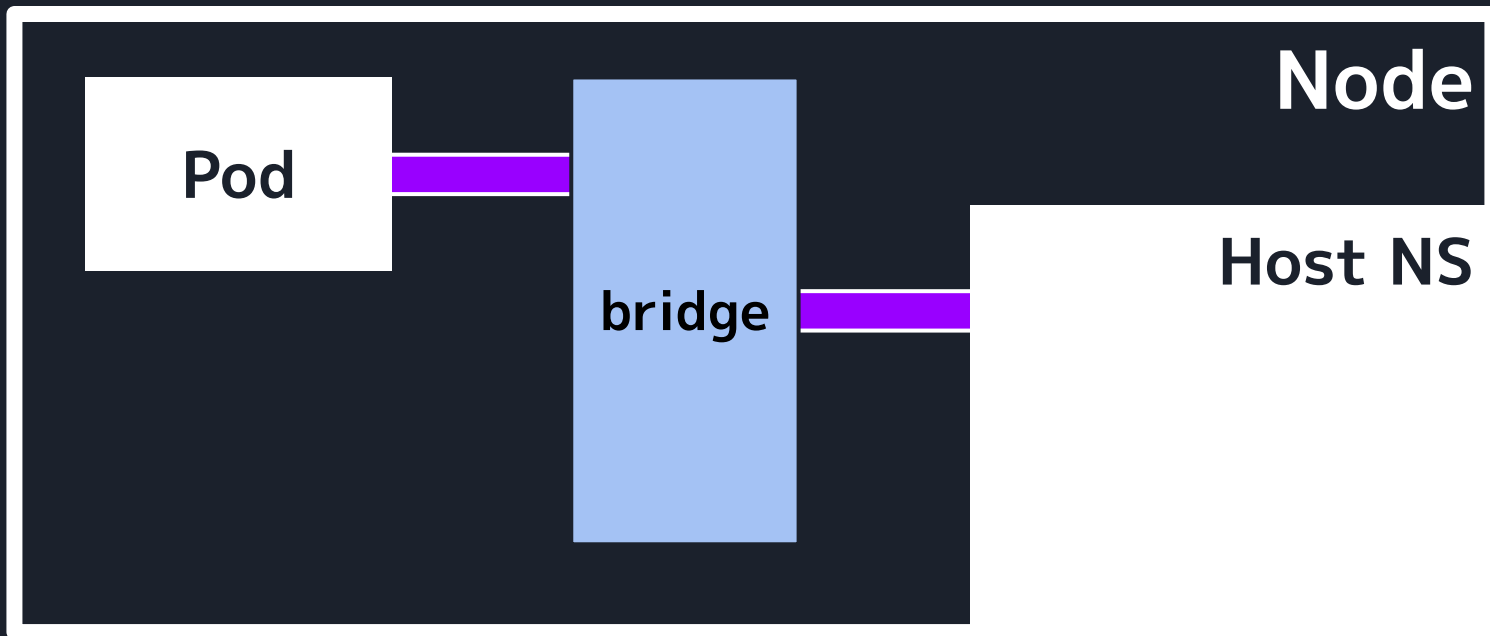
# Flannel - 同じノード内

Pod作成時



# Flannel - 同じノード内

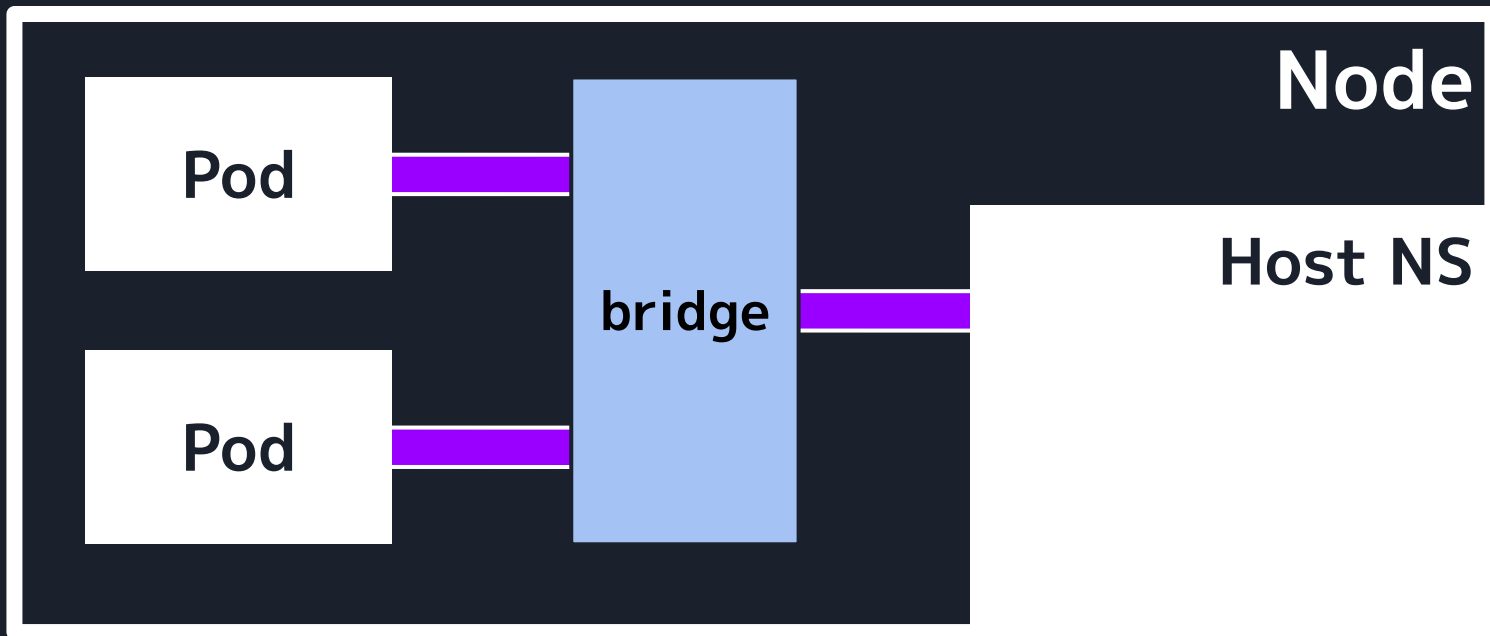
Pod作成時: Linux Bridgeに接続





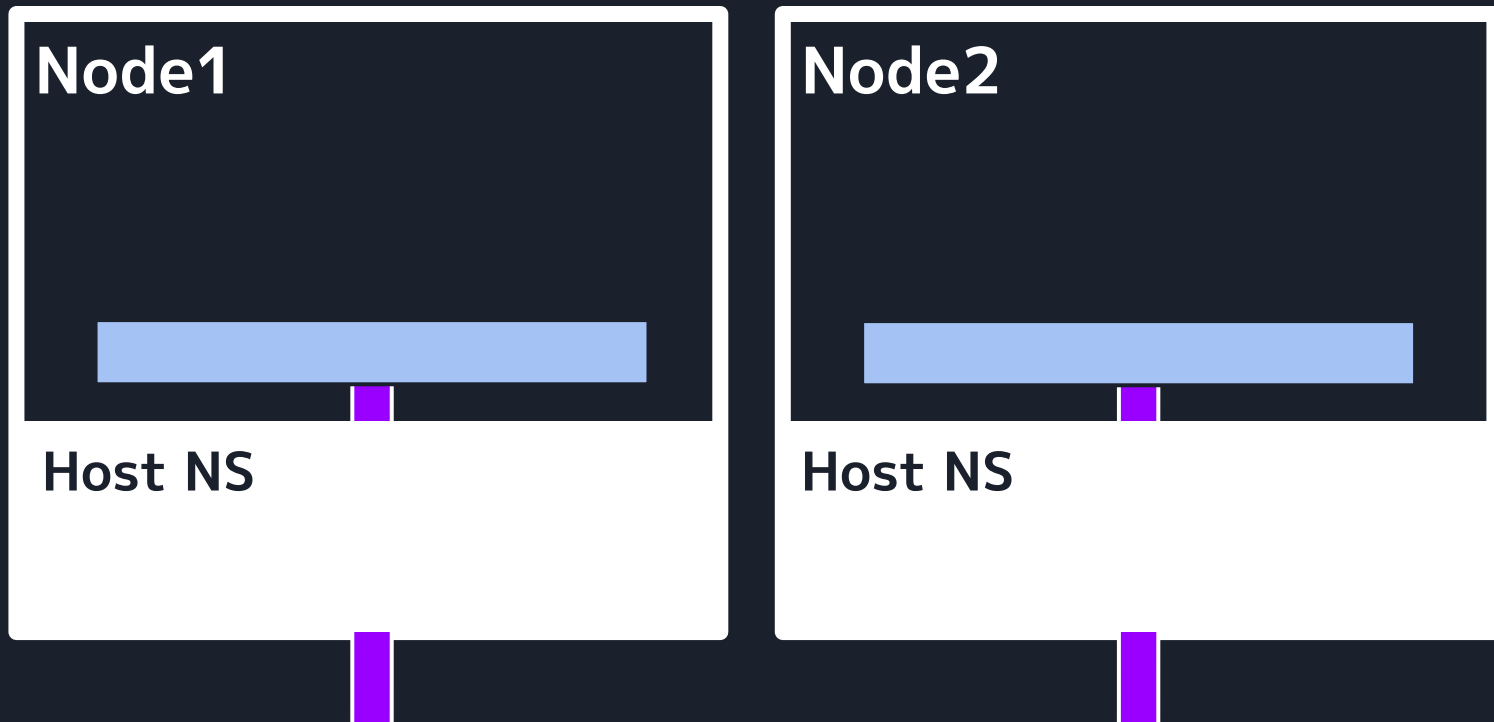
# Flannel - 同じノード内

通信時: 全てのノード内Podは**L2で接続**され、疎通ができる



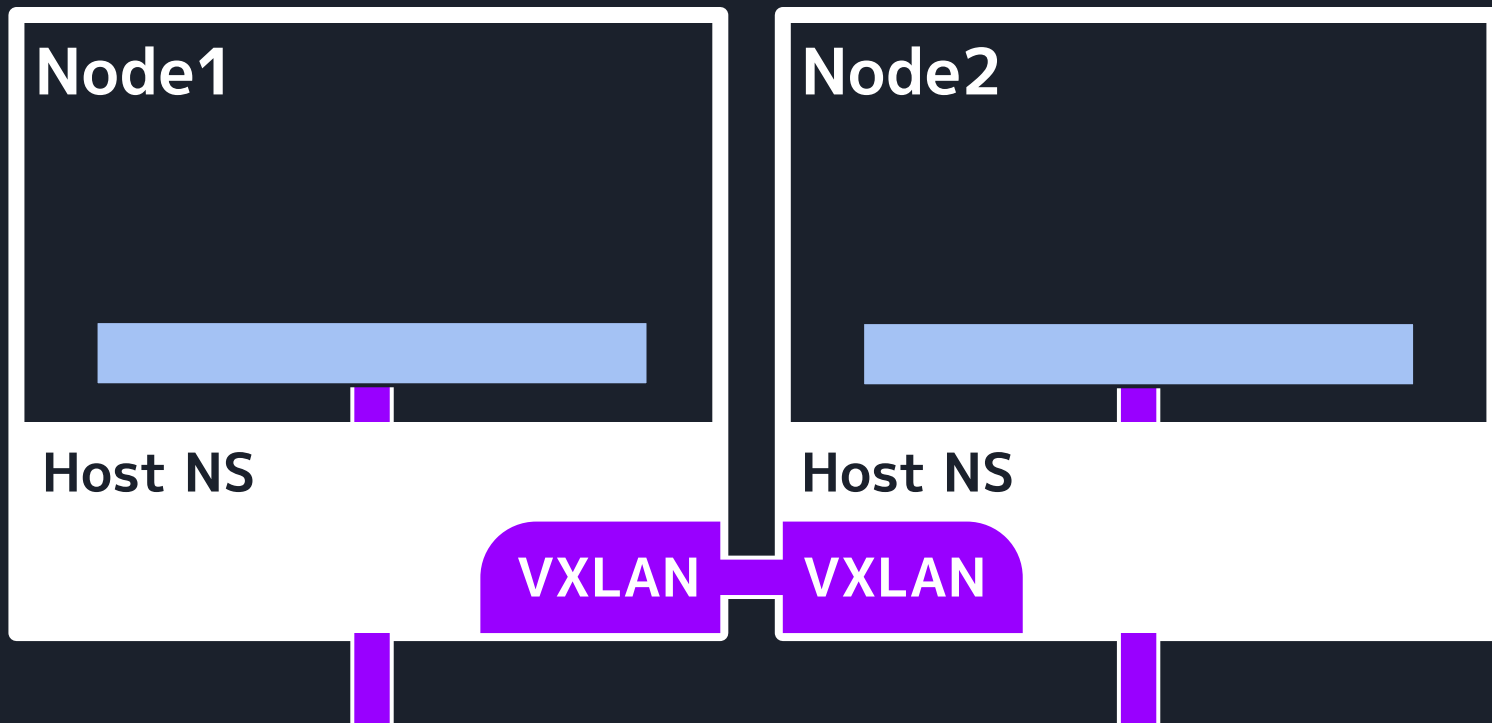
# Flannel - 異なるノード間

下準備: 同じノード内の下準備が終わった段階



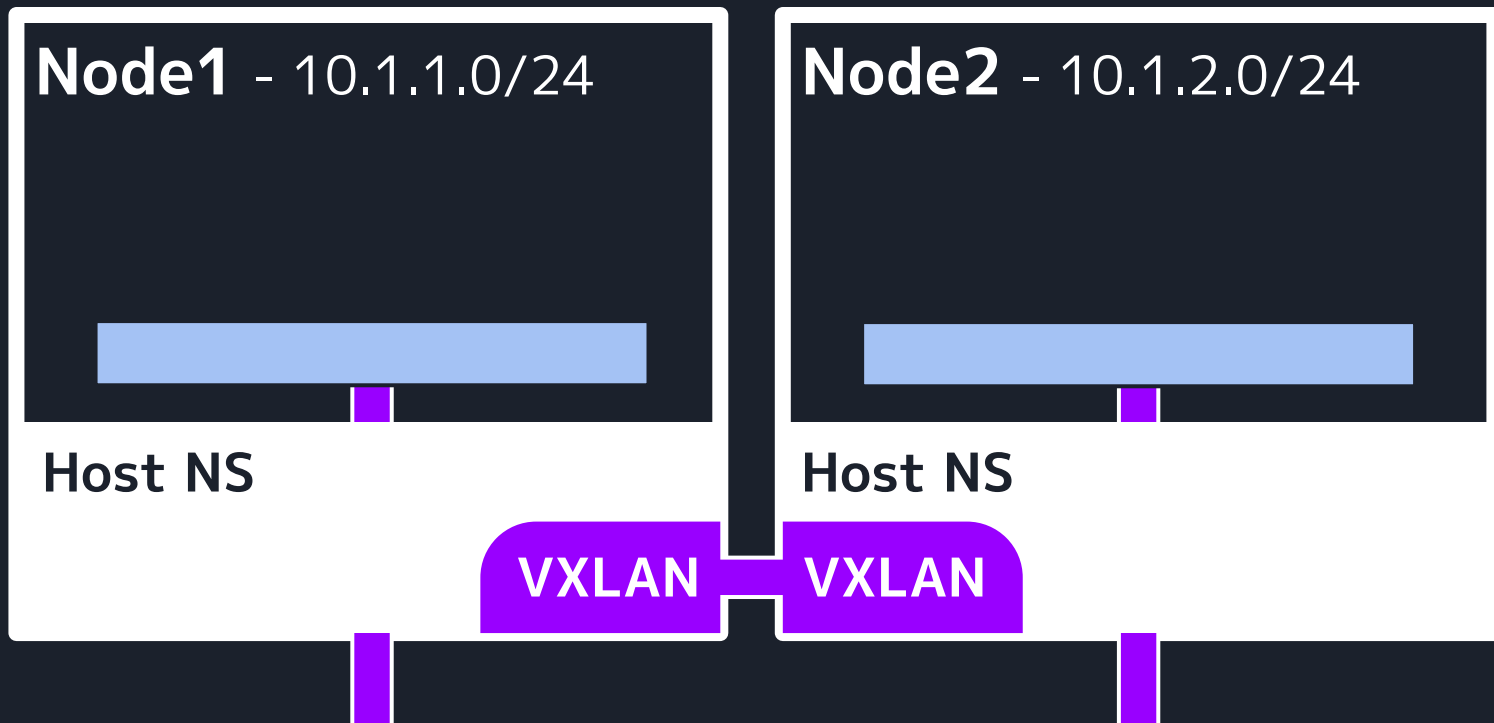
# Flannel - 異なるノード間

下準備: 各ノードでVTEPを作り、VXLANで接続



# Flannel - 異なるノード間

下準備: 各ノードに固有のサブネットを割り当てる

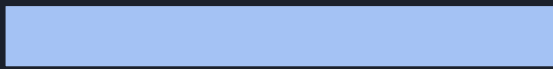


# Flannel - 異なるノード間

下準備: 各ノードに固有のサブネット

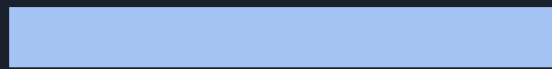
Pod全体のIPレンジは  
**10.1.0.0/16**とする

**Node1** - 10.1.1.0/24



Host NS

**Node2** - 10.1.2.0/24



Host NS

VXLAN

VXLAN



# Flannel - 異なるノード間

下準備: 各ノードに固有のサブネットを割り当てる

**Node1** - 10.1.1.0/24

**Node2** - 10.1.2.0/24

このサブネットで  
どのノードにいるPodなのか  
**判別**する

Host NS

Host NS

VXLAN

VXLAN



# Flannel - 異なるノード間

下準備: ルーティングテーブルにエントリを追加

**Node1** - 10.1.1.0/24

destination	gateway
10.1.1.0/24	Bridge
10.1.0.0/16	VTEP

**Node2** - 10.1.2.0/24

destination	gateway
10.1.2.0/24	Bridge
10.1.0.0/16	VTEP

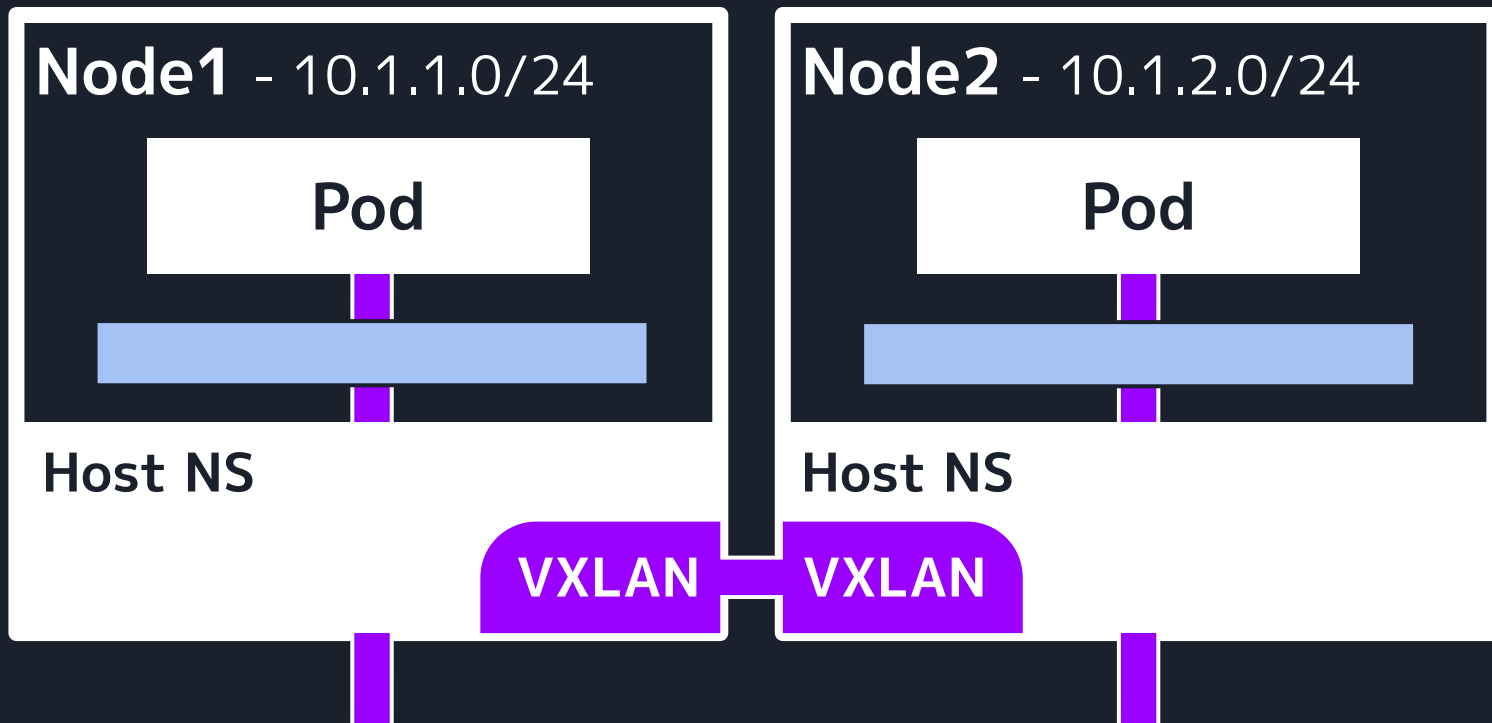
Host 1

VXLAN

VXLAN

# Flannel - 異なるノード間

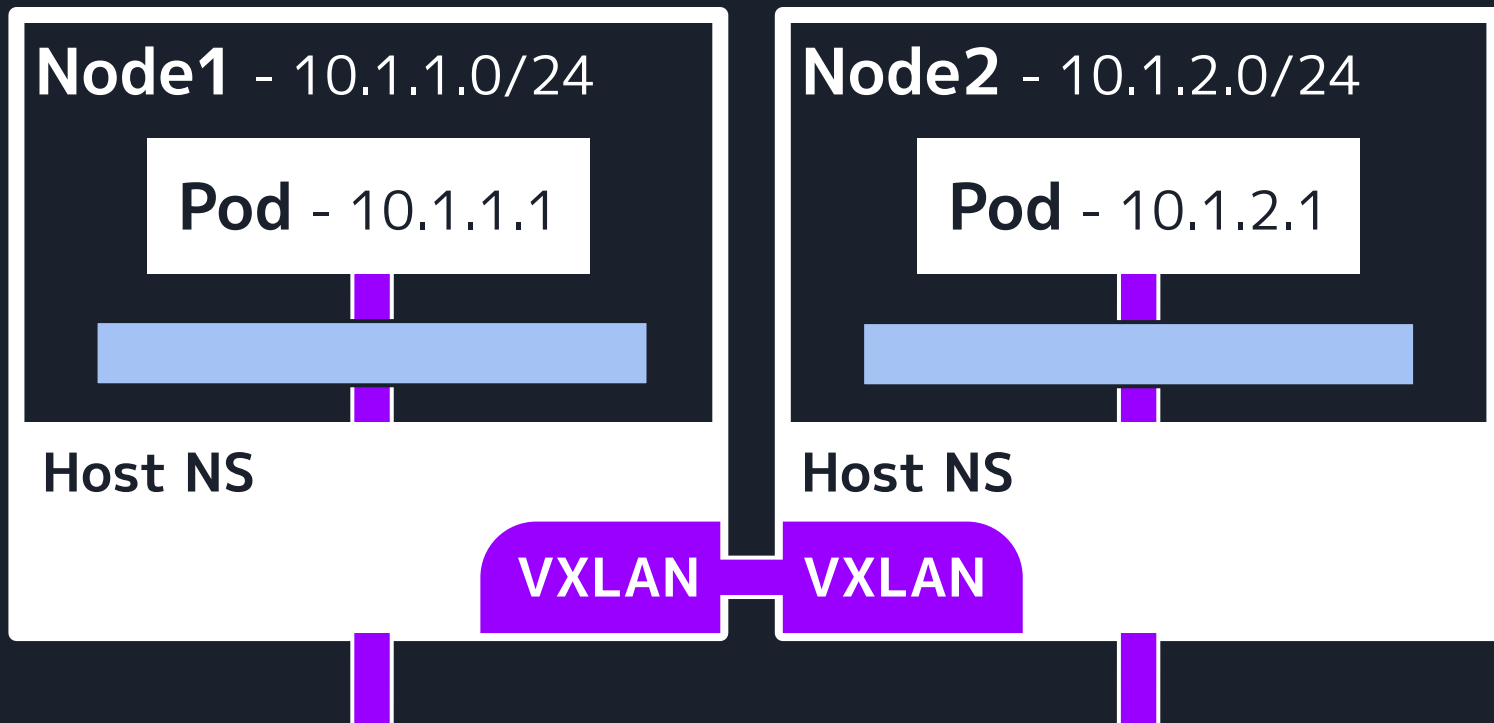
Pod作成時





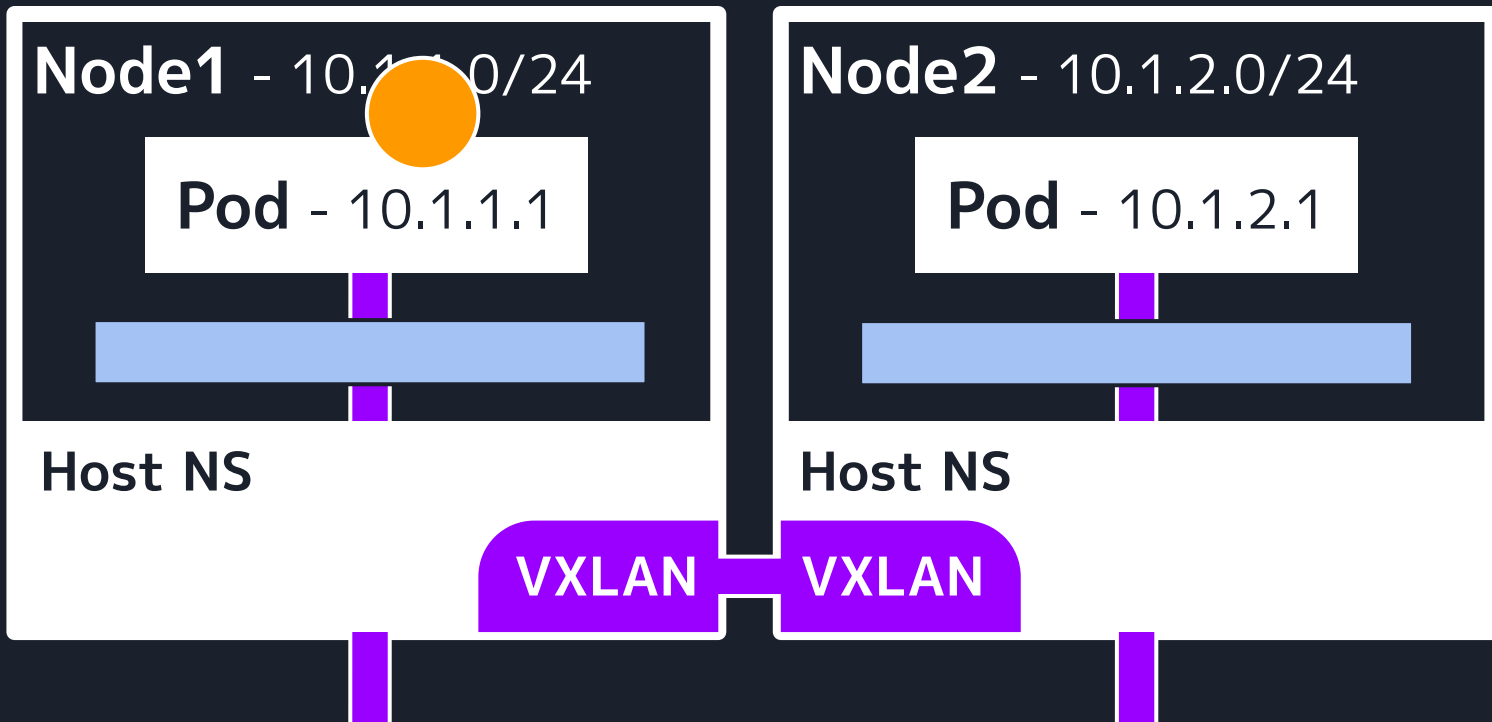
# Flannel - 異なるノード間

Pod作成時: ノードのIPレンジからアドレスを割り当てる



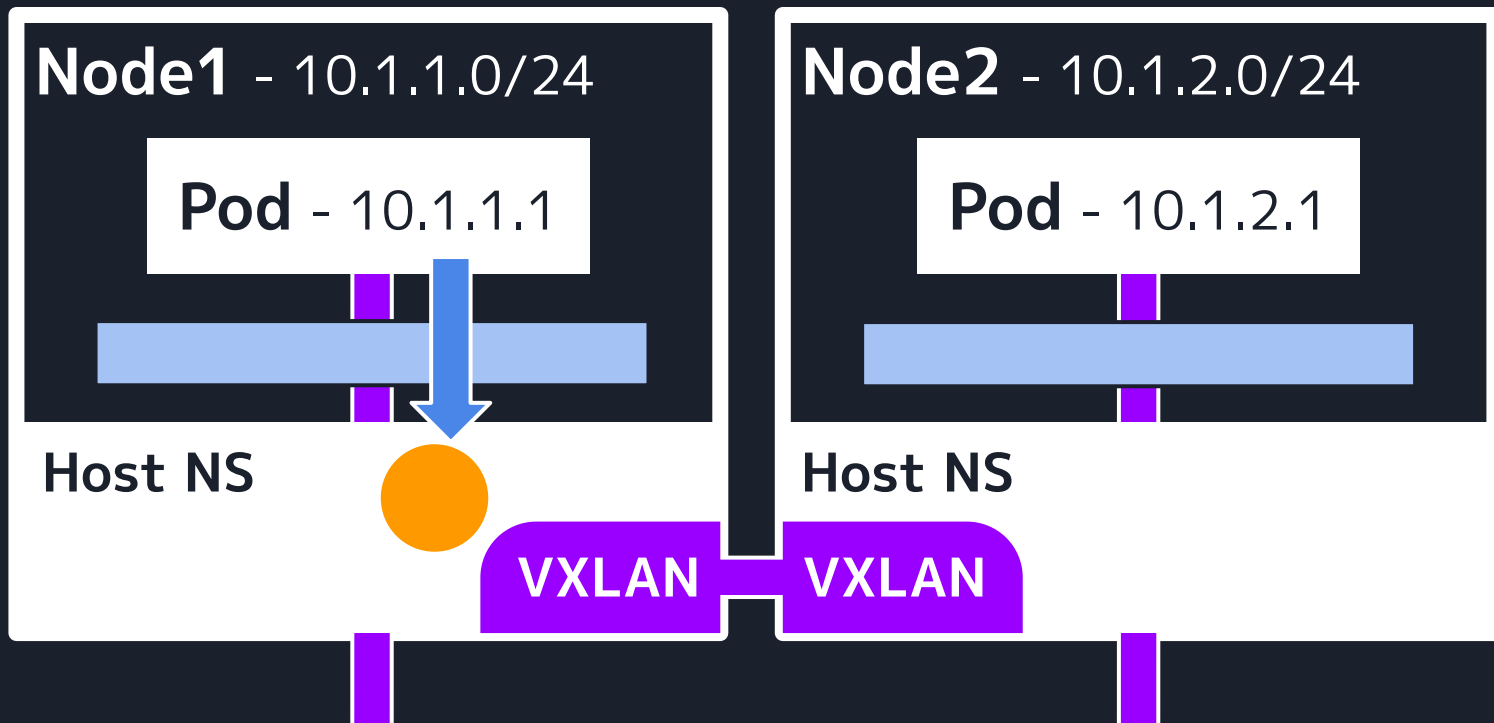
# Flannel - 異なるノード間

通信時: 10.1.1.1から10.1.2.1へ送る時



# Flannel - 異なるノード間

通信時: デフォルトゲートウェイとなっているHost NSへ



# Flannel - 異なるノード間

通信時: ルーティングテーブルを確認

**Node1** - 10.1.1.0/24

destination	gateway
10.1.1.0/24	Bridge
<u>10.1.0.0/16</u>	<u>VTEP</u>

**Node2** - 10.1.2.0/24

**Pod** - 10.1.2.1

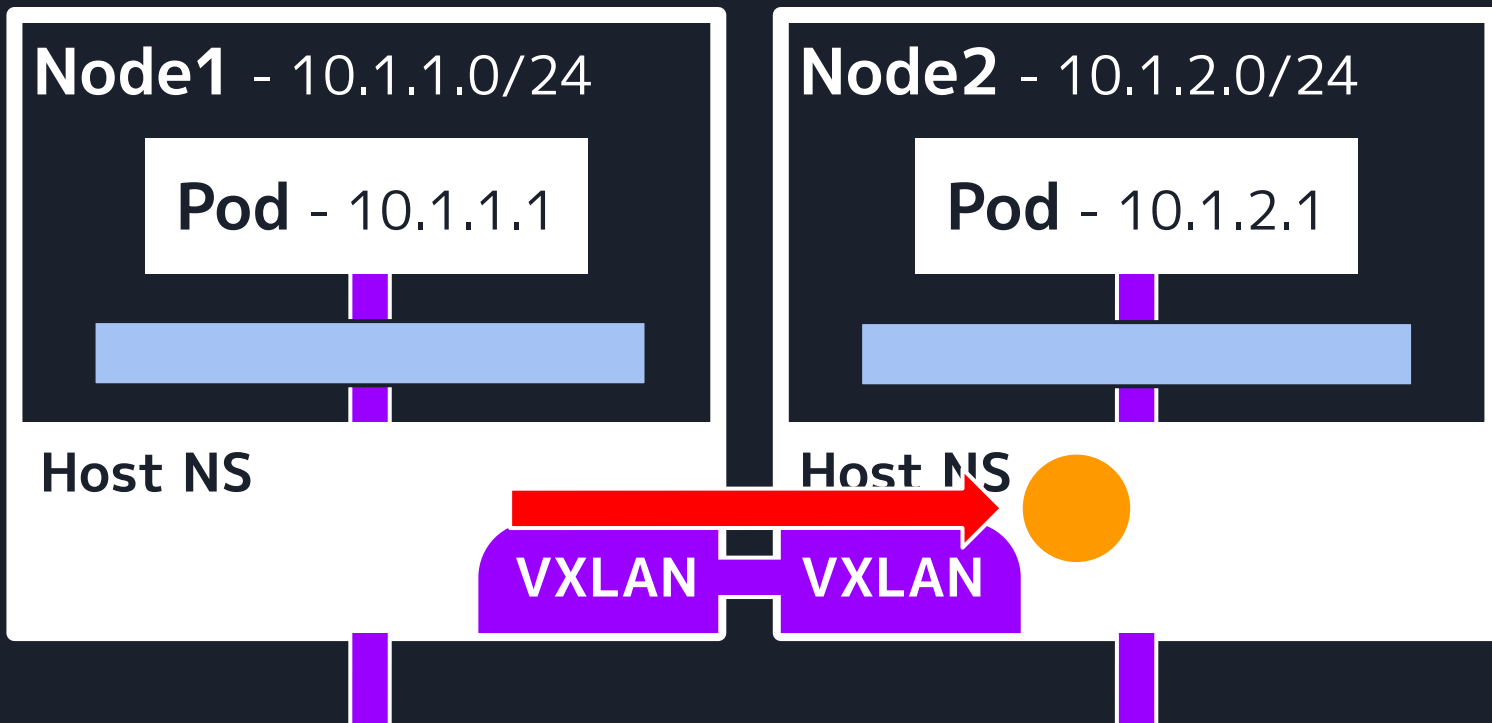
Host NS

VXLAN

VXLAN

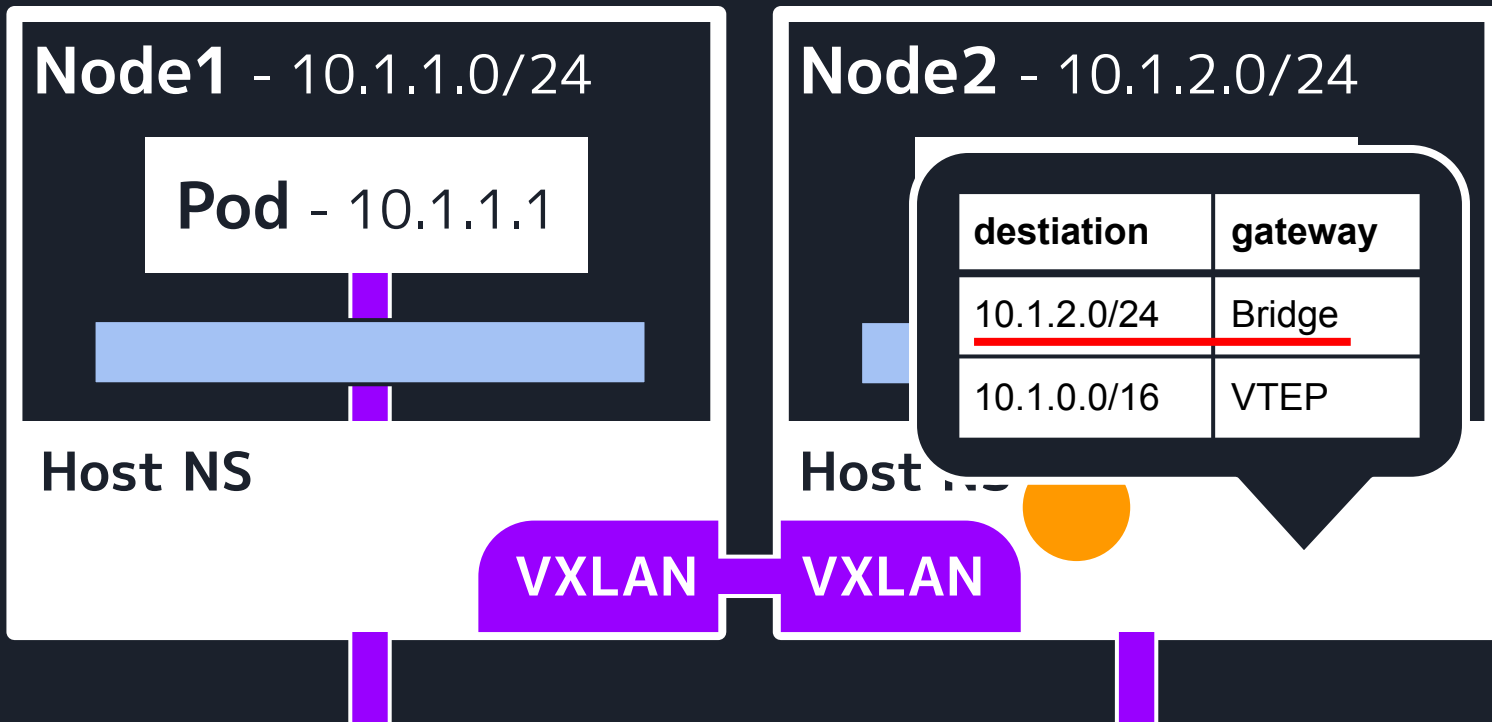
# Flannel - 異なるノード間

通信時: gatewayの**VTEP**から**VXLAN**経由でNode2へ送信



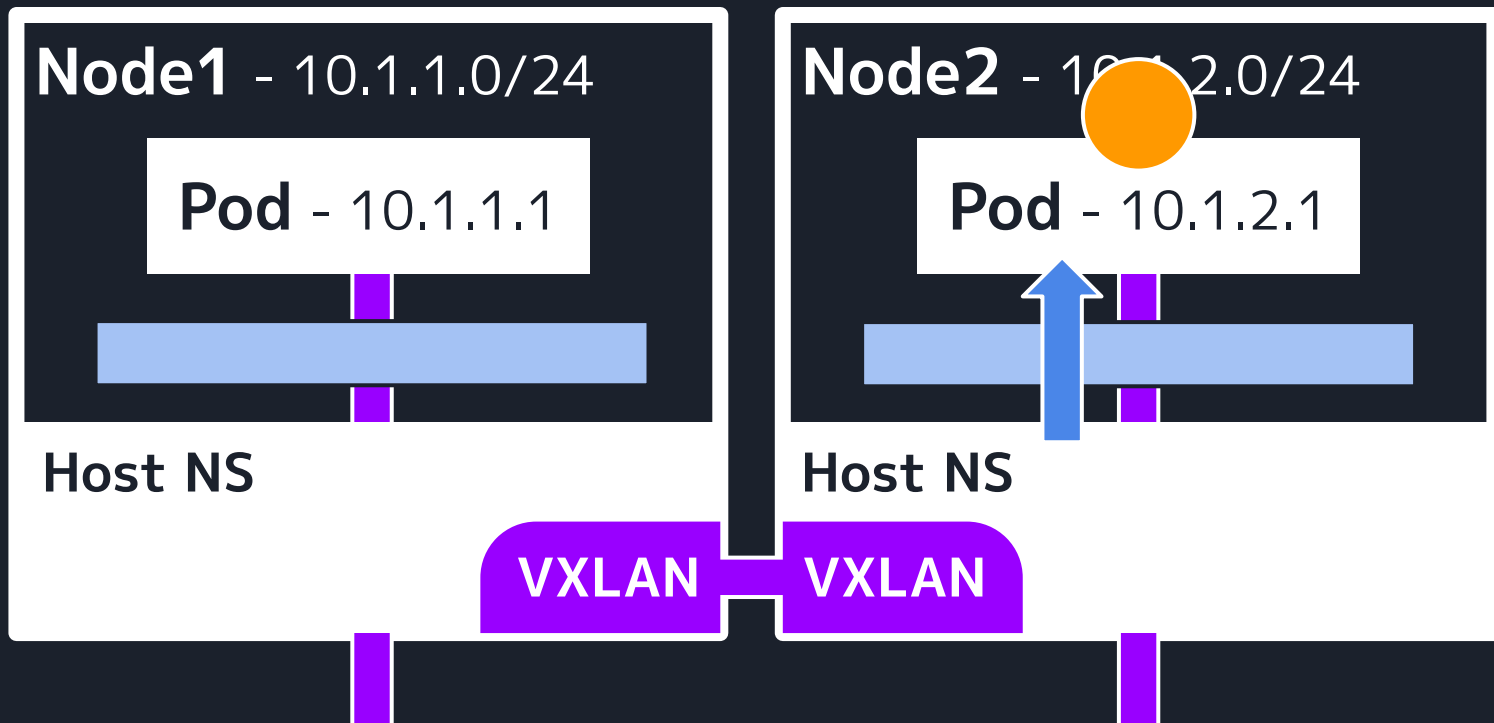
# Flannel - 異なるノード間

通信時: ルーティングテーブルを確認



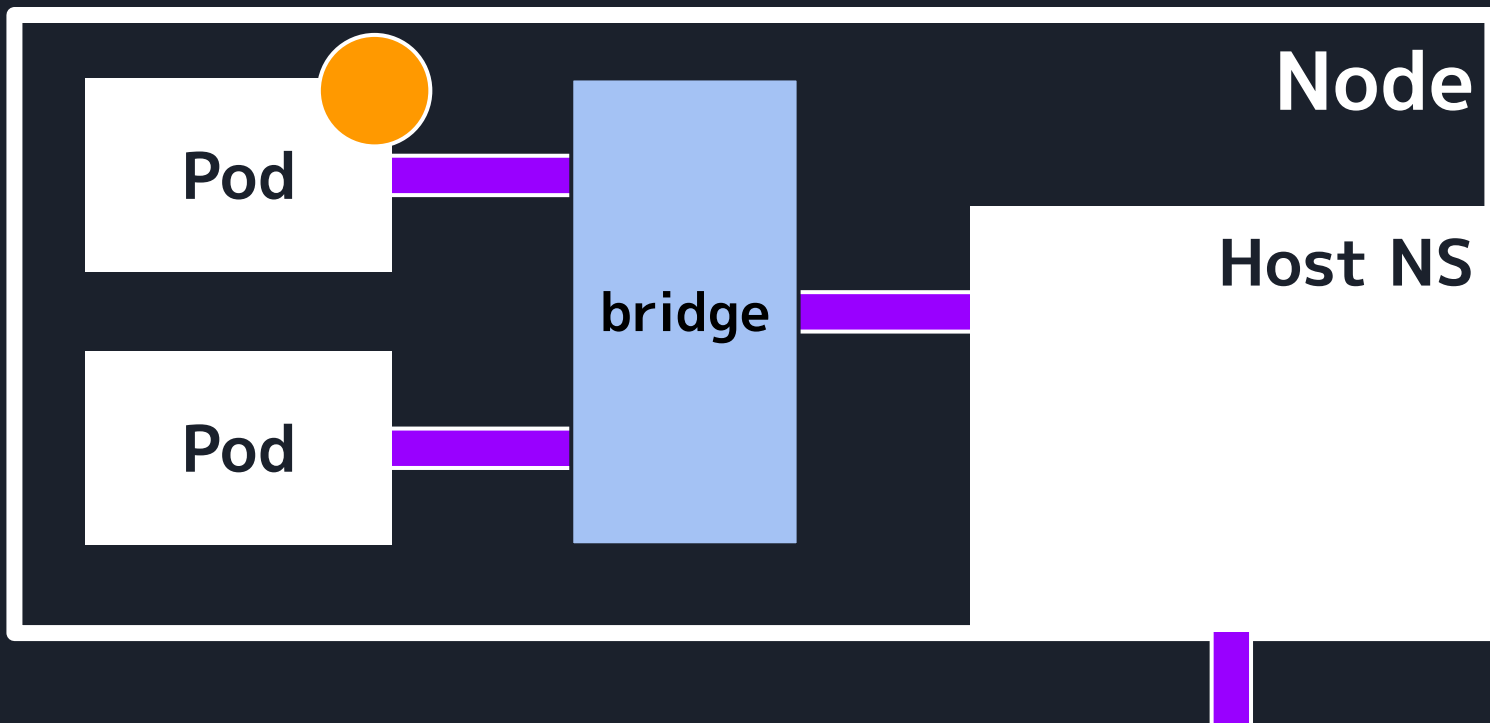
# Flannel - 異なるノード間

通信時: gatewayのBridgeを経由し、宛先に到達



# Flannel - Podから外部ネットワーク

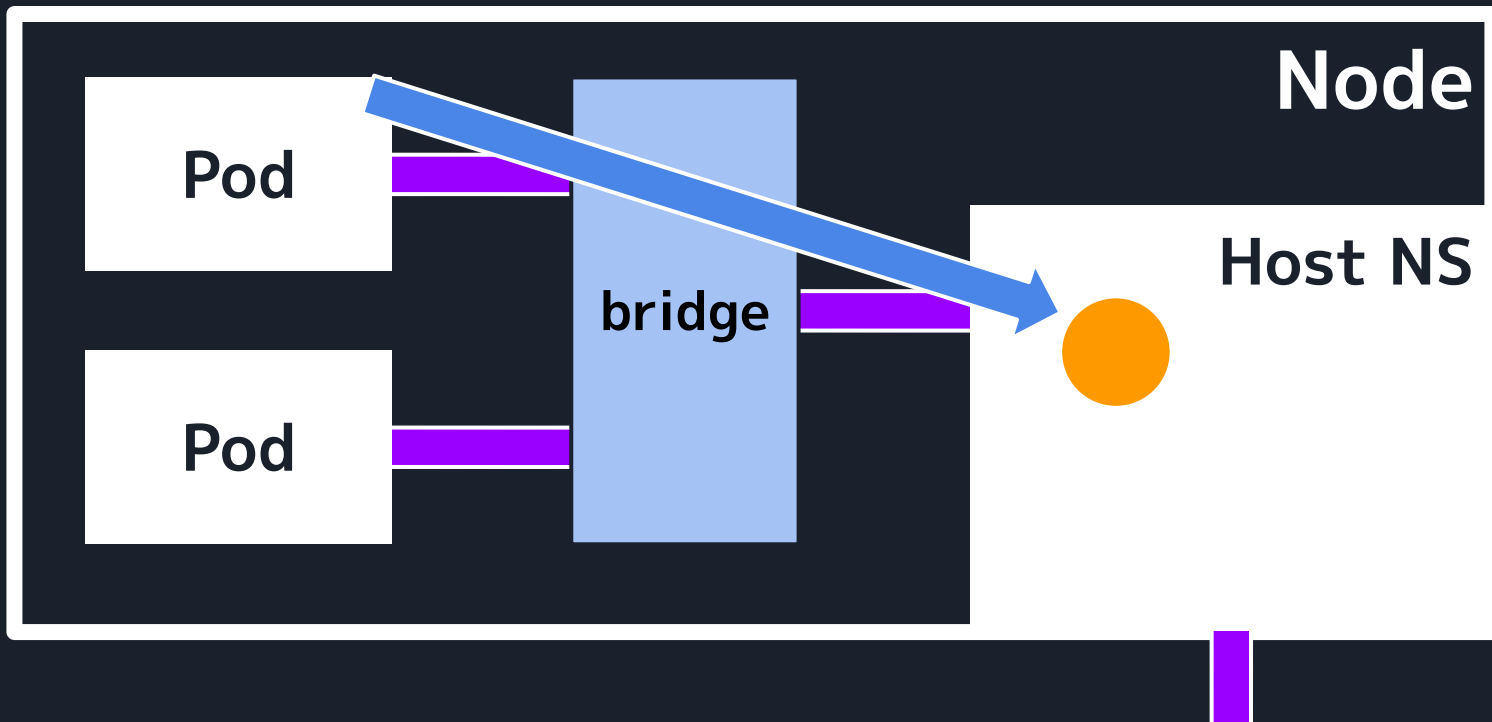
NetfilterでIPマスカレードをして外部に出す





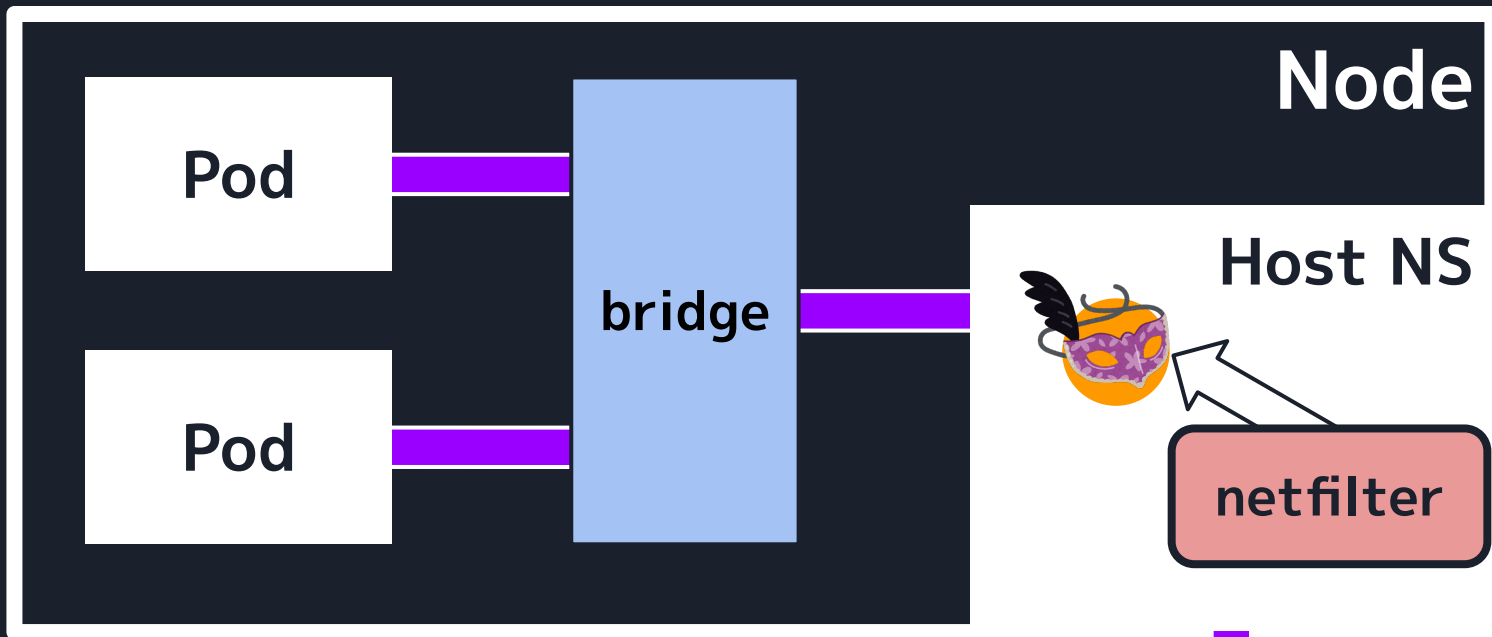
# Flannel - Podから外部ネットワーク

NetfilterでIPマスカレードをして外部に出す



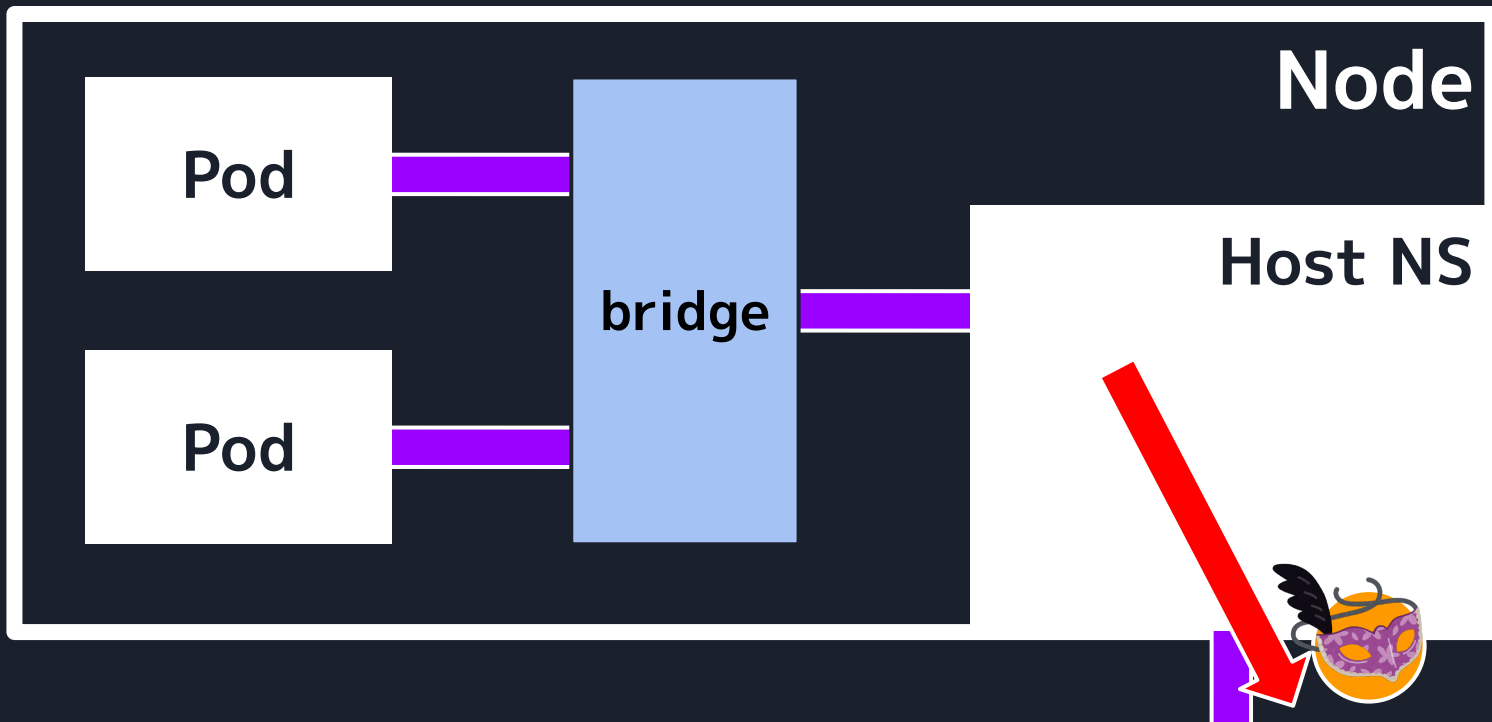
# Flannel - Podから外部ネットワーク

NetfilterでIPマスカレードをして外部に出す



# Flannel - Podから外部ネットワーク

NetfilterでIPマスカレードをして外部に出す





# Flannel - 特徴まとめ

- **同じノード内でのL3疎通**
  - **Linux Bridge**で全てのPodをL2接続する
- **異なるノード間のL3疎通**
  - **VXLAN**を用いて別ノードにパケットを丸ごと輸送する  
**オーバーレイネットワーク**
    - 余談) Wireguardなど**VPN**で輸送するモードも
- **サブネット**でノードを見分ける
- **netfilter**でIPマスカレードする



# Calico - 同じノード内

下準備: Flannelと同じくノードにサブネットを割り当て

**Node1** - 10.1.1.0/24

**Host NS**



# Calico - 同じノード内

Pod作成時

**Node1** - 10.1.1.0/24

**Pod**

10.1.1.1

**Host NS**

# Calico - 同じノード内

Pod作成時: **veth**でHost Namespaceと**直接接続**

**Node1** - 10.1.1.0/24

**Pod**  
10.1.1.1

NIC1

**Host NS**



# Calico - 同じノード内

Pod作成時: ルーティングテーブルにエントリを追加

**Node1** - 10.1.1.0/24

**Pod**  
10.1.1.1

NIC1

**Host NS**

destiation	gateway
10.1.1.1/32	NIC1



# Calico - 同じノード内

通信時: スタティックルーティングでL3疎通ができる

**Node1** - 10.1.1.0/24

**Pod**  
10.1.1.1

**Pod**  
10.1.1.2

NIC1

NIC2

**Host NS**

destiation	gateway
10.1.1.1/32	NIC1
10.1.1.2/32	NIC2



# Calico - 異なるノード間

下準備: 同じノード内の下準備が終わった段階

**Node1** - 10.1.1.0/24

Host NS

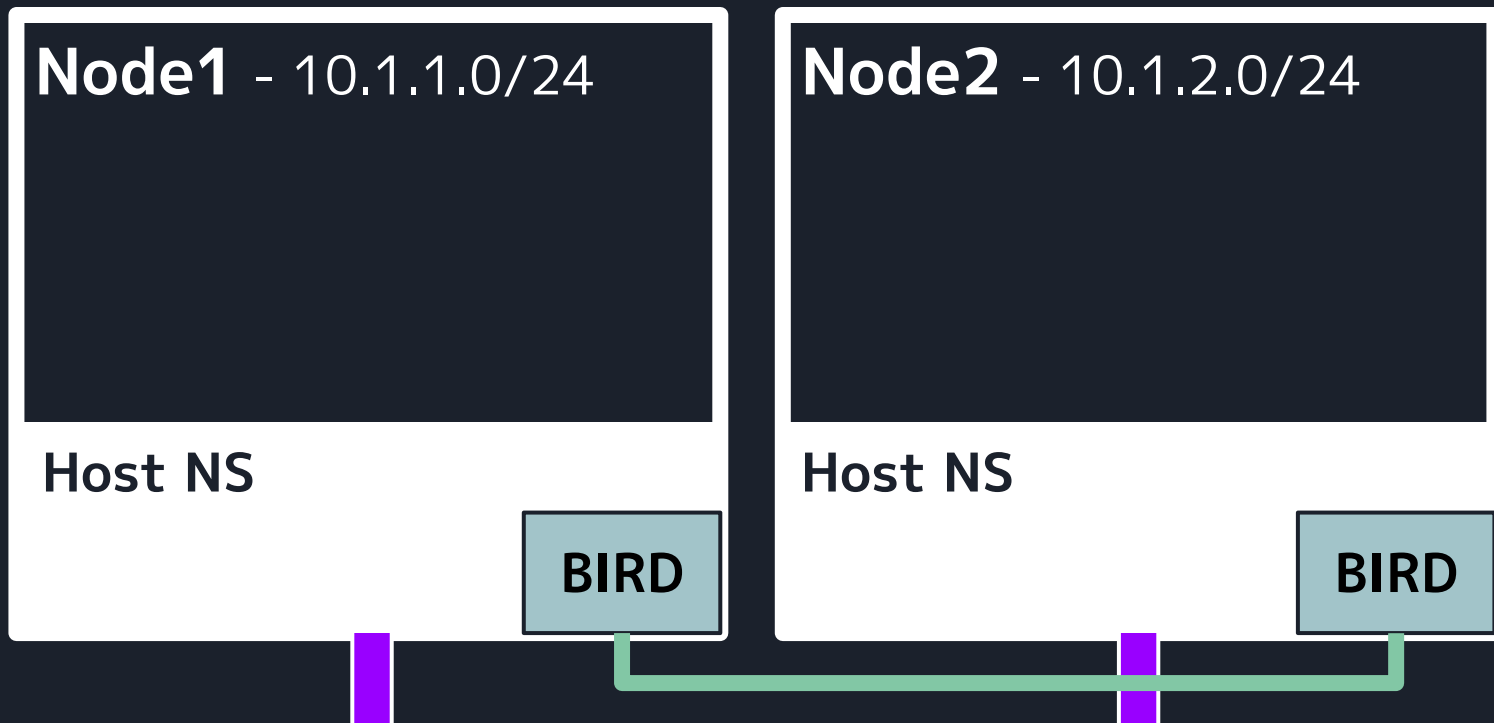
**Node2** - 10.1.2.0/24

Host NS



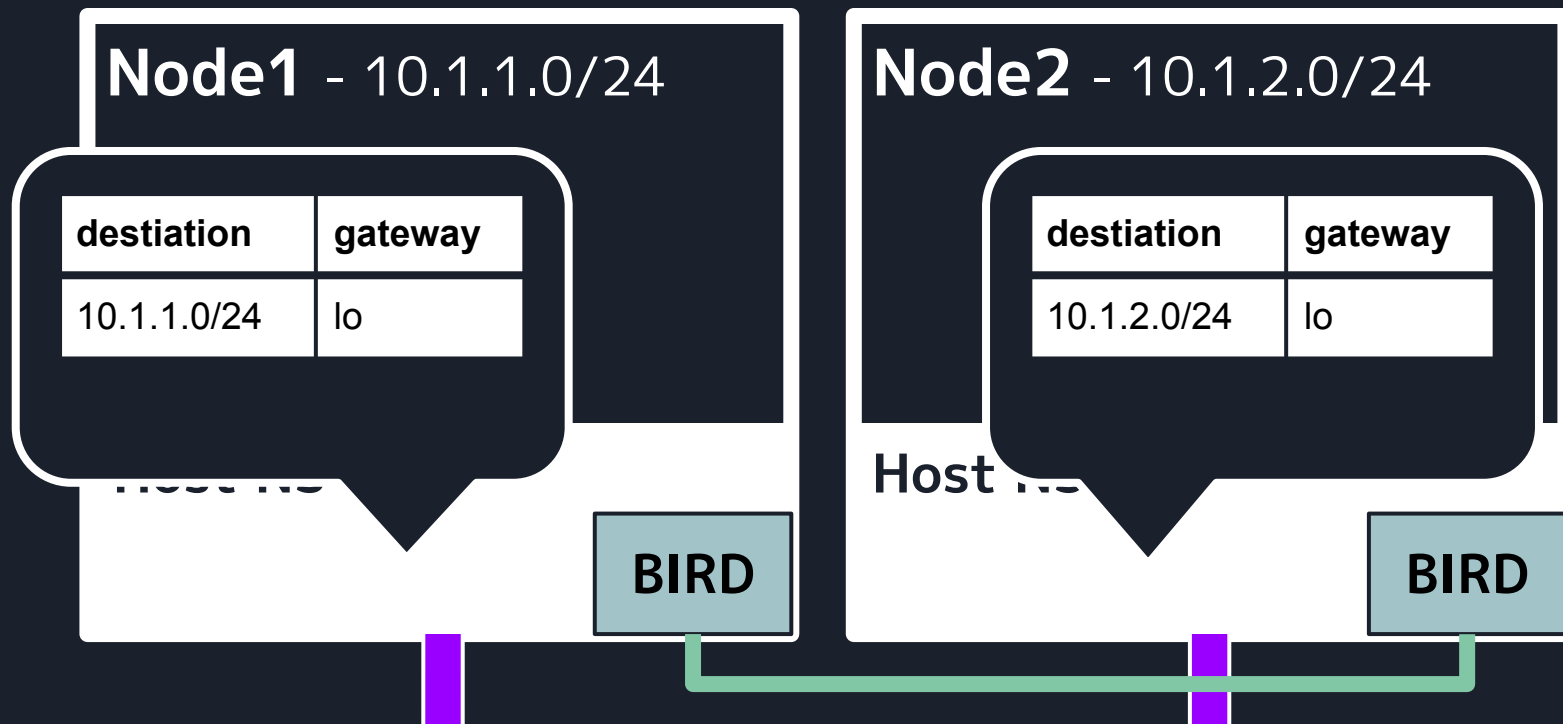
# Calico - 異なるノード間

下準備: BIRD (ルーターソフトウェア) で**BGPピア**を張る



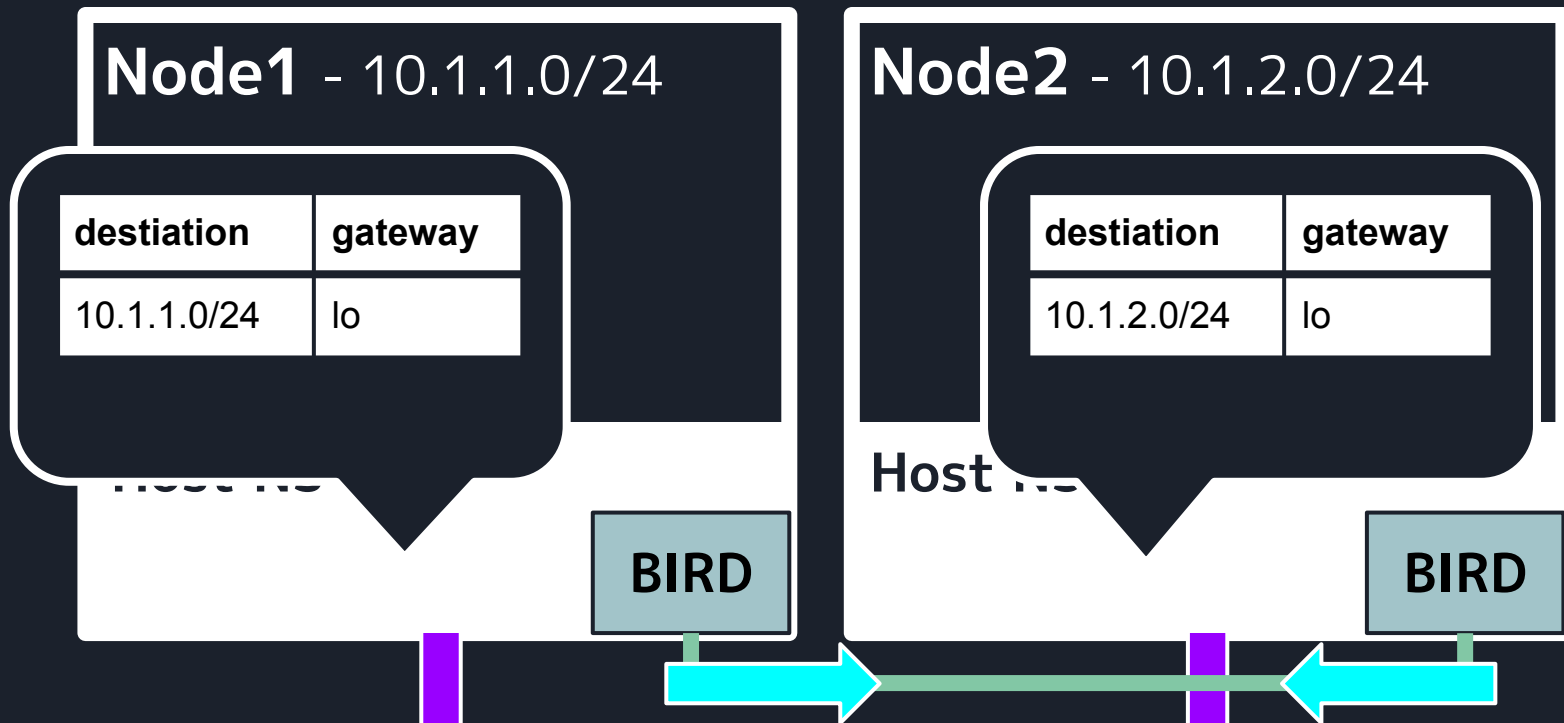
# Calico - 異なるノード間

下準備: 自分のノードのPodサブネットを互いに**広報**する



# Calico - 異なるノード間

下準備: 自分のノードのPodサブネットを互いに**広報**する



# Calico - 異なるノード間

下準備: 自分のノードのPodサブネットを互いに広報する

**Node1** - 10.1.1.0/24

destination	gateway
10.1.1.0/24	lo
<u>10.1.2.0/24</u>	<u>Node2</u>

BIRD

**Node2** - 10.1.2.0/24

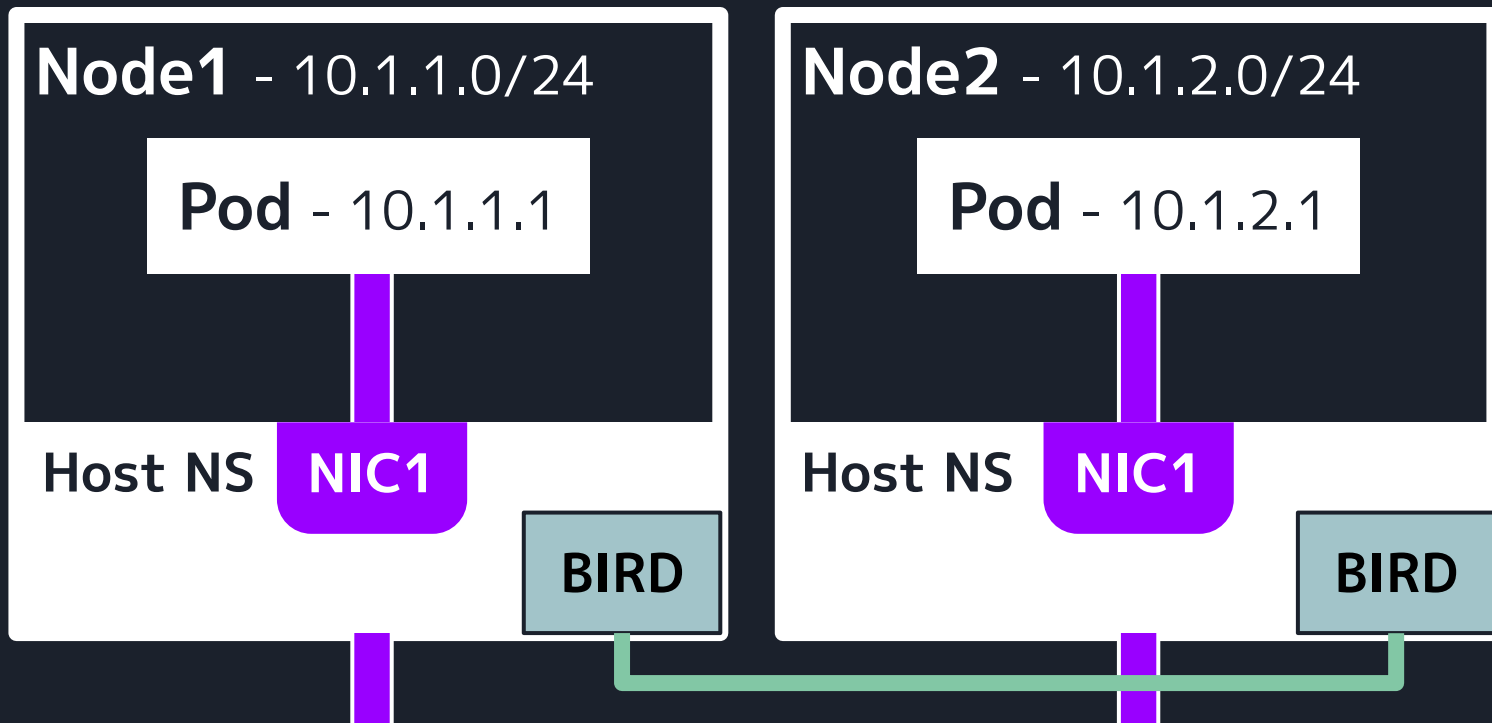
destination	gateway
10.1.2.0/24	lo
<u>10.1.1.0/24</u>	<u>Node1</u>

Host ...

BIRD

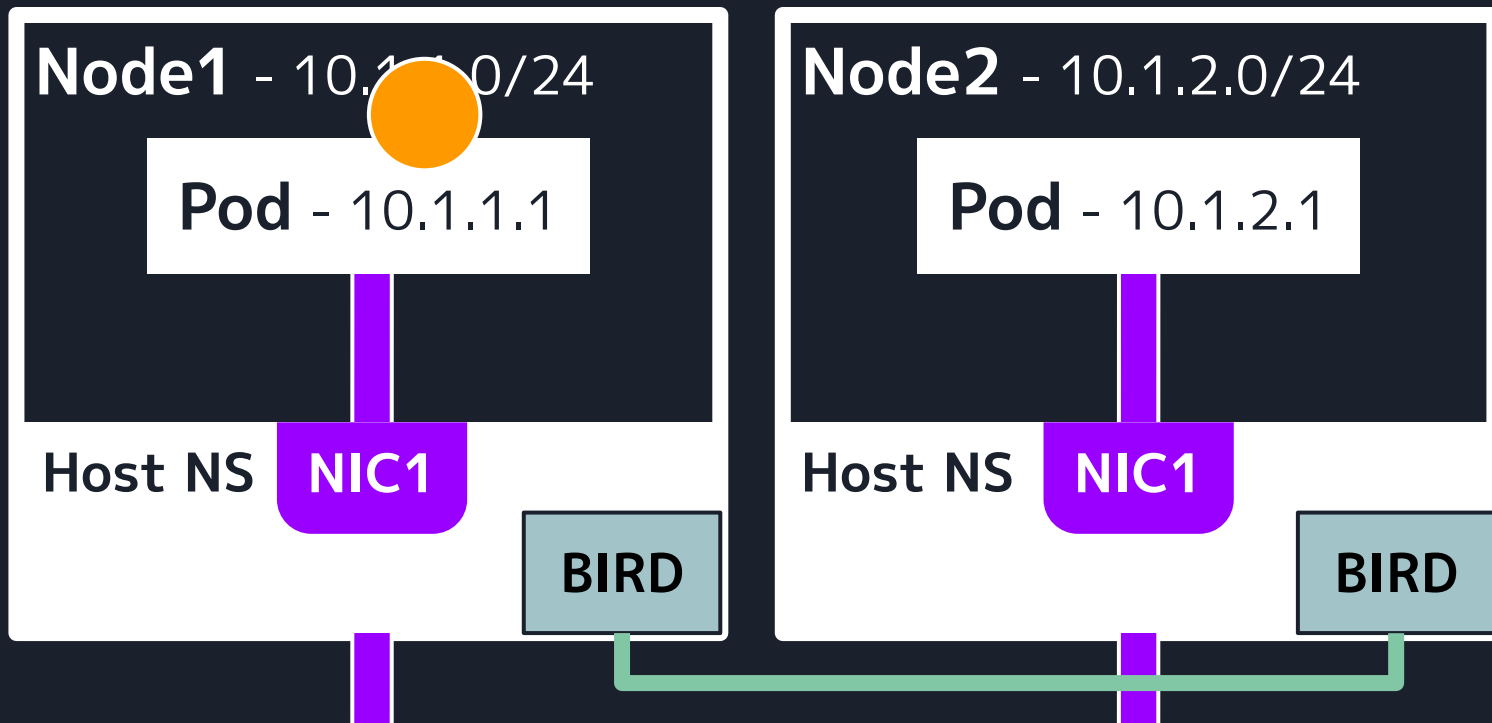
# Calico - 異なるノード間

Pod作成時: 追加操作は特になし



# Calico - 異なるノード間

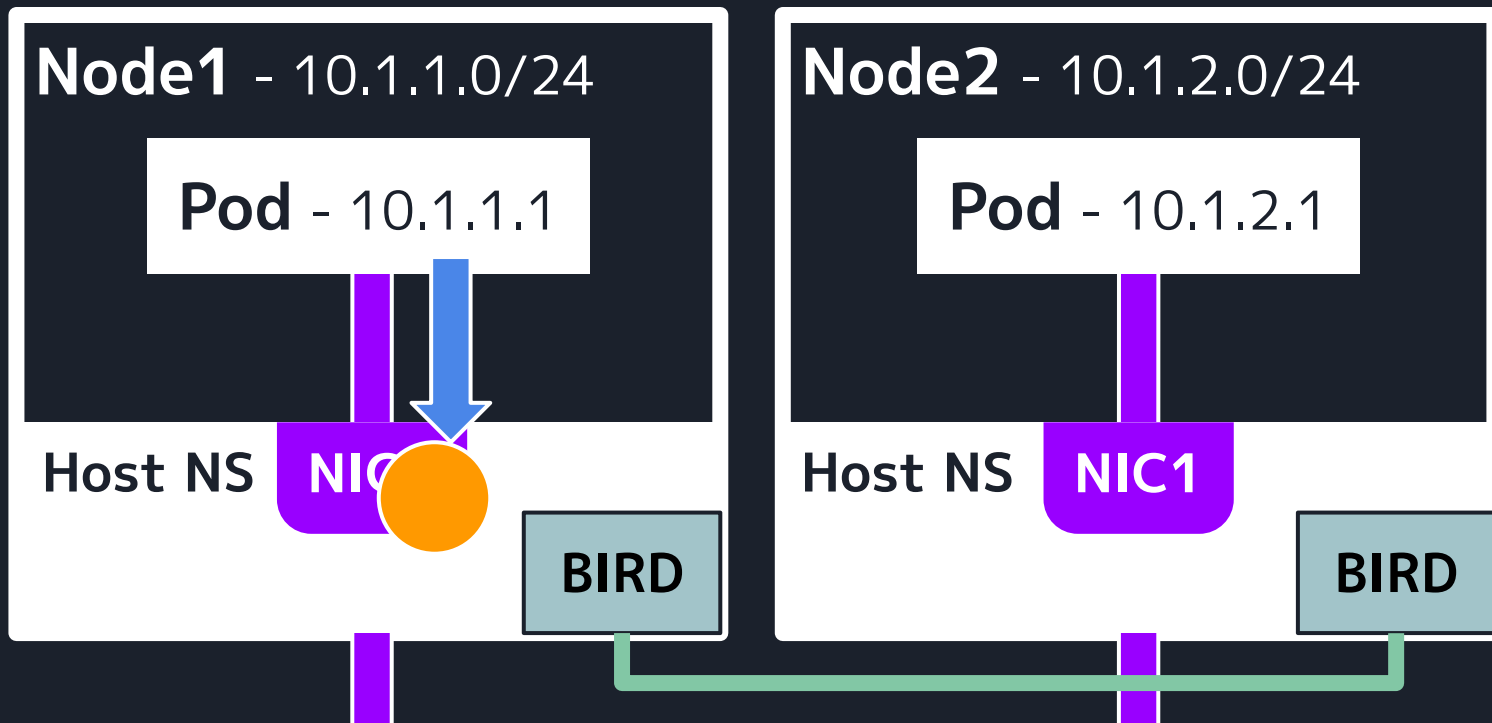
通信時: 10.1.1.1から10.1.2.1へ送る時





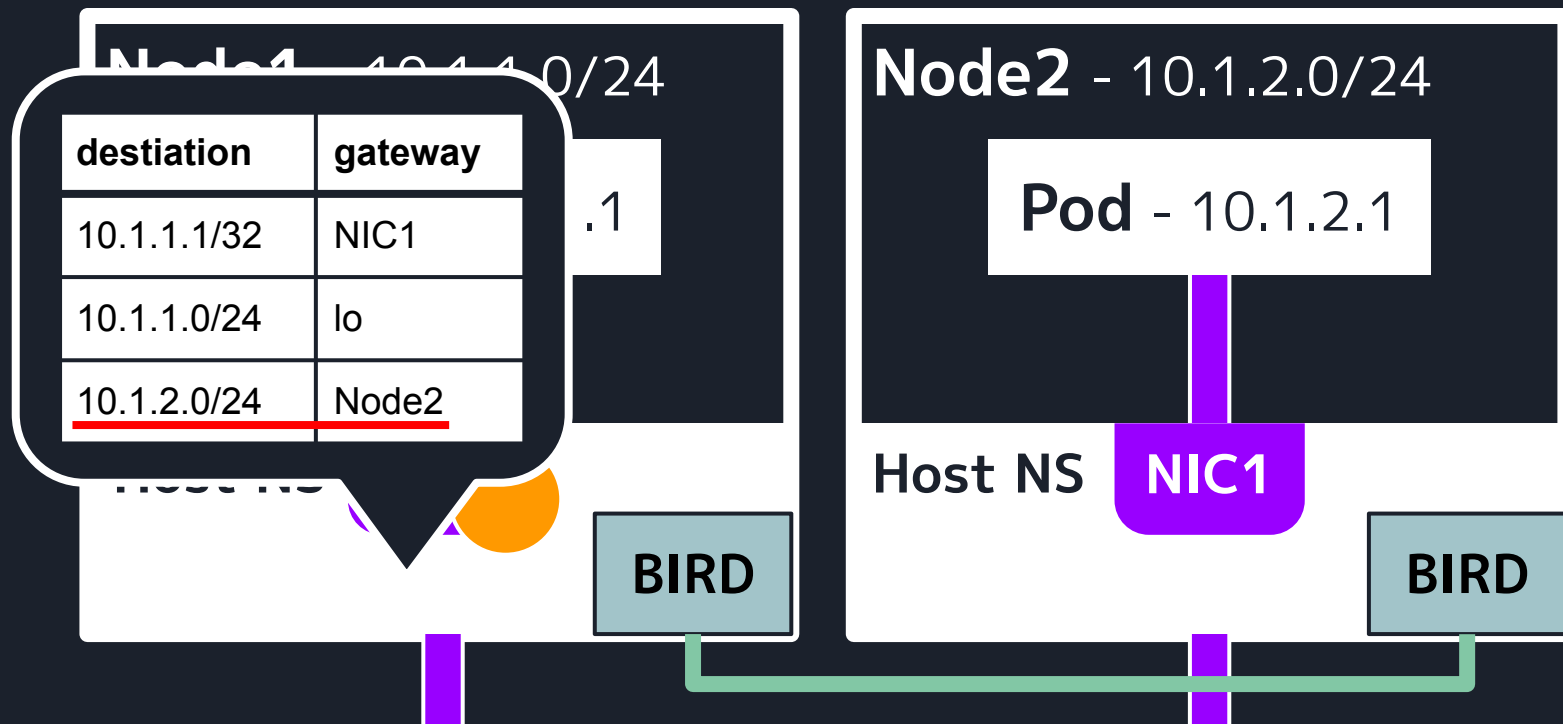
# Calico - 異なるノード間

通信時: デフォルトゲートウェイとなっているHost NSへ



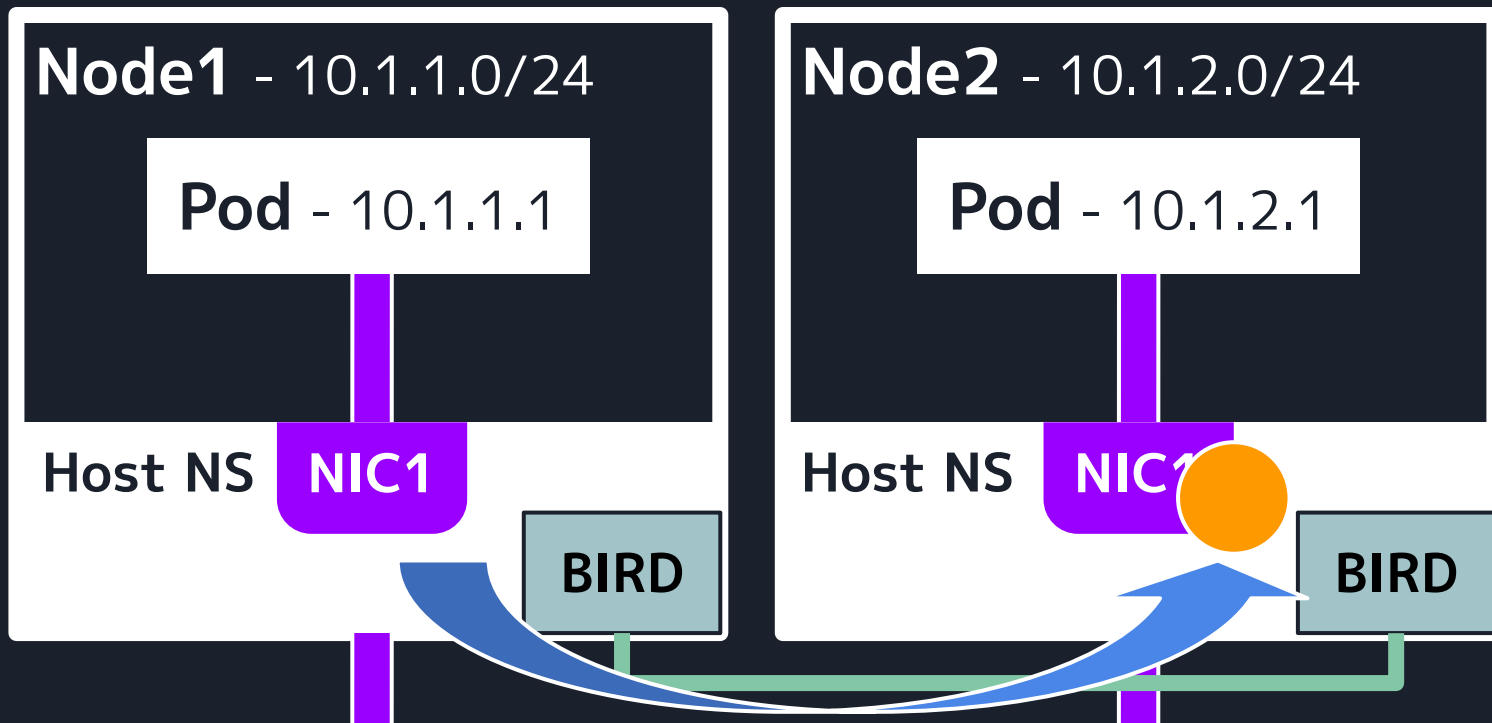
# Calico - 異なるノード間

通信時: ルーティングテーブルを確認



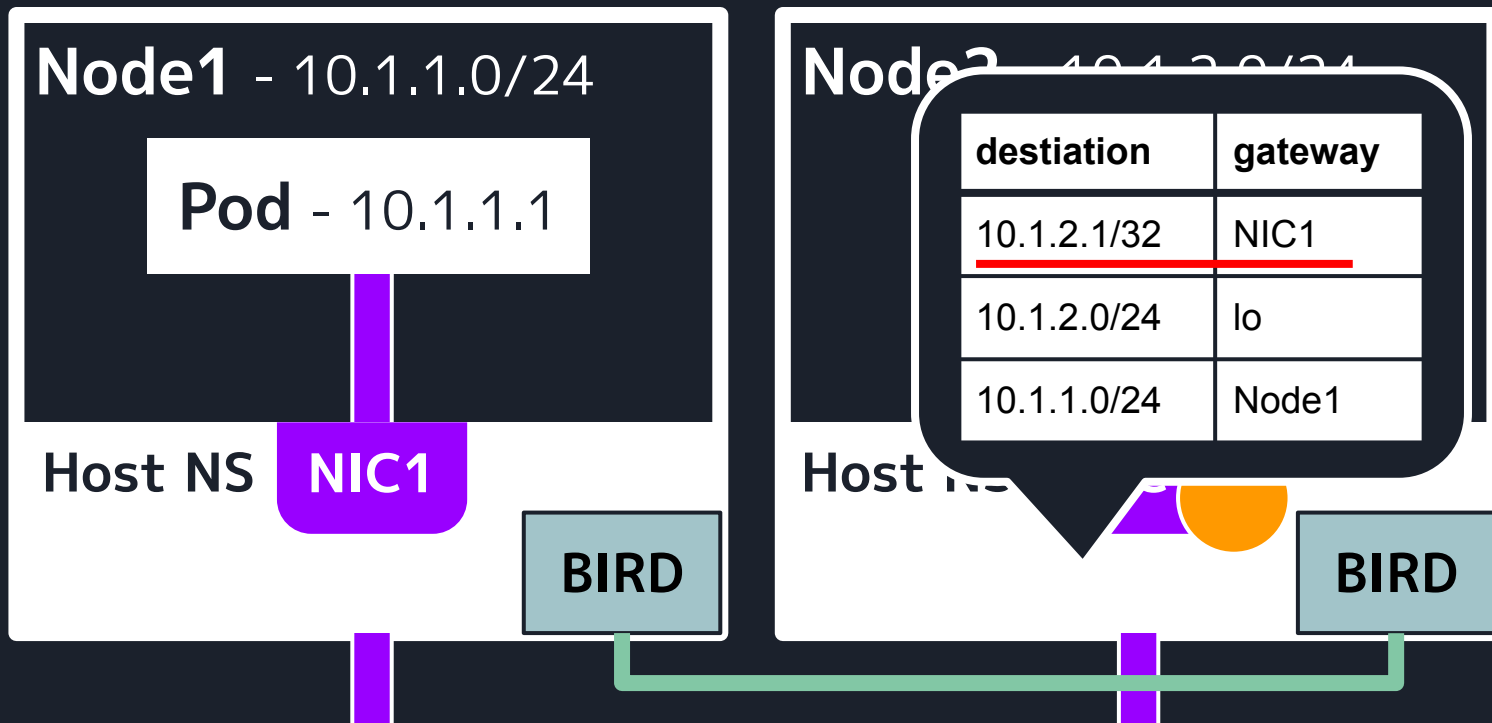
# Calico - 異なるノード間

通信時: アンダーレイから直接、gatewayのNode2へ送信



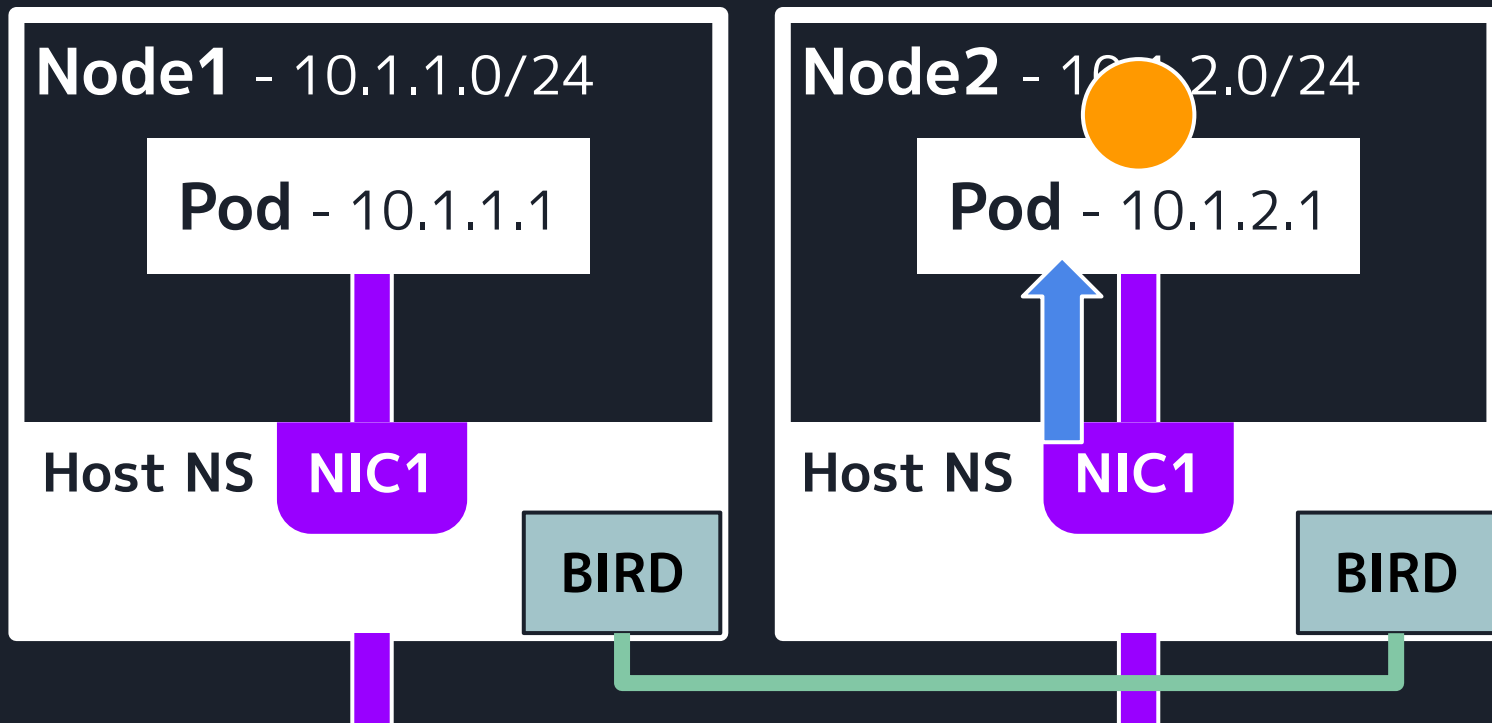
# Calico - 異なるノード間

通信時: ルーティングテーブルを確認



# Calico - 異なるノード間

通信時: gatewayの**NIC1**から宛先に到達



# Calico - Podから外部ネットワーク

Flannel同様、Netfilterで**IPマスカレード**をして外部に出す

**Node1** - 10.1.1.0/24

**Pod**  
10.1.1.1

**Pod**  
10.1.1.2

NIC1

NIC2



**Host NS**

netfilter



# Calico - 特徴まとめ

- 同じノード内でのL3疎通
  - スタティックルーティングする
- 異なるノード間のL3疎通
  - BGPで経路を広報、直接相手ノードに送るPure L3
    - オーバーレイよりも高効率
  - 【制約】全てのノードが同じL2セグメントに属している必要がある
- それ以外はFlannel同様



# Cilium - 同じノード内

下準備: 前二つと同じくノードにサブネットを割り当て

**Node1** - 10.1.1.0/24

**Host NS**





# Cilium - 同じノード内

Pod作成時

**Node1** - 10.1.1.0/24

**Pod**

10.1.1.1

**Host NS**

# Cilium - 同じノード内

Pod作成時: **veth**でHost Namespaceと**直接接続**

**Node1** - 10.1.1.0/24

**Pod**  
10.1.1.1

NIC1

**Host NS**



# Cilium - 同じノード内

Pod作成時: Host側のNICに**eBPF**がアタッチされる

**Node1** - 10.1.1.0/24

**Pod**  
10.1.1.1

NIC1



**Host NS**

# Cilium - 同じノード内

Pod作成時: 別Podも同様

**Node1** - 10.1.1.0/24

**Pod**  
10.1.1.1

**Pod**  
10.1.1.2

NIC1



Host NS

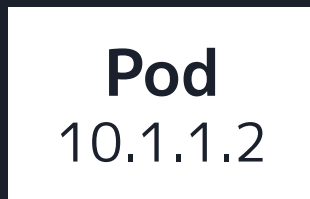
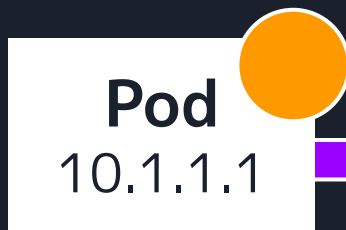
NIC2



# Cilium - 同じノード内

通信時: 10.1.1.1から10.1.1.2へ送る時

**Node1** - 10.1.1.0/24



Host NS

# Cilium - 同じノード内

通信時: デフォルトゲートウェイとなっているHost NSへ

**Node1** - 10.1.1.0/24

**Pod**  
10.1.1.1

**Pod**  
10.1.1.2



NIC1



**Host NS**

NIC2



# Cilium - 同じノード内

通信時: パケットがNICに着くとeBPFプログラムが起動

**Node1** - 10.1.1.0/24

**Pod**  
10.1.1.1

**Pod**  
10.1.1.2



NIC1



Host NS

NIC2



# Cilium - 同じノード内

通信時: 宛先MACを宛先PodのMACに書き換える

**Node1** - 10.1.1.0/24

**Pod**  
10.1.1.1

**Pod**  
10.1.1.2



**Host NS**





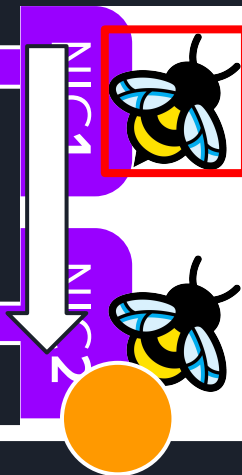
# Cilium - 同じノード内

通信時: eBPFの機能でNIC2に直接転送

**Node1** - 10.1.1.0/24

**Pod**  
10.1.1.1

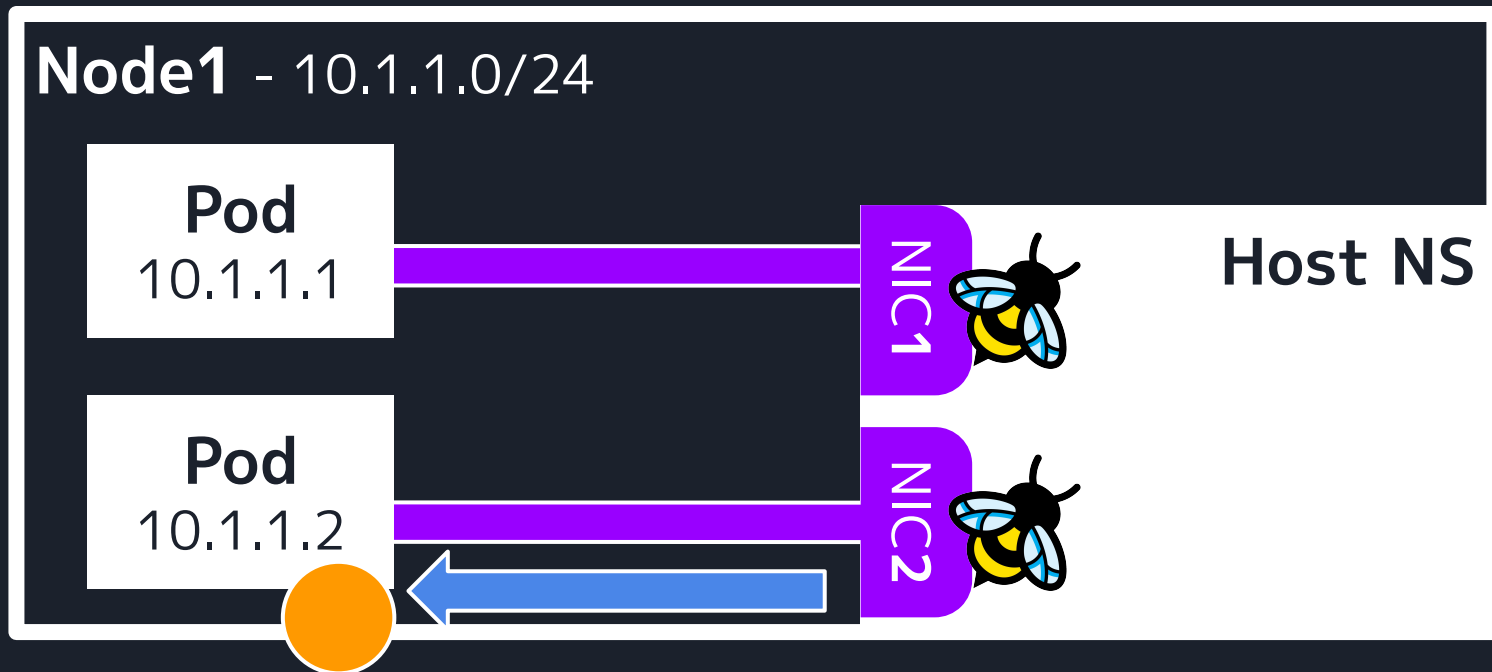
**Pod**  
10.1.1.2



**Host NS**

# Cilium - 同じノード内

通信時: あとはそのまま宛先Podへ





## Cilium - 異なるノード間

- Ciliumは、同じノードのPod宛でなければ、eBPFプログラムを抜けて**Linux側に処理を投げる**
- Ciliumはあらかじめ、異なるノード間のパケット転送のために以下のいずれかを用意する
  - **Overlay**モード: Flannelと同じ**VXLAN**を使う転送
  - **Native**モード: Calicoと同じ**BGP広報**を使う転送
- すなわち、Ciliumの異なるノード間疎通は**FlannelかCalicoと全く同じ仕組み**で行われる

# Cilium - Podから外部ネットワーク

eBPFで**IPマスカレード**をして外部に出す

**Node1** - 10.1.1.0/24

**Pod**  
10.1.1.1

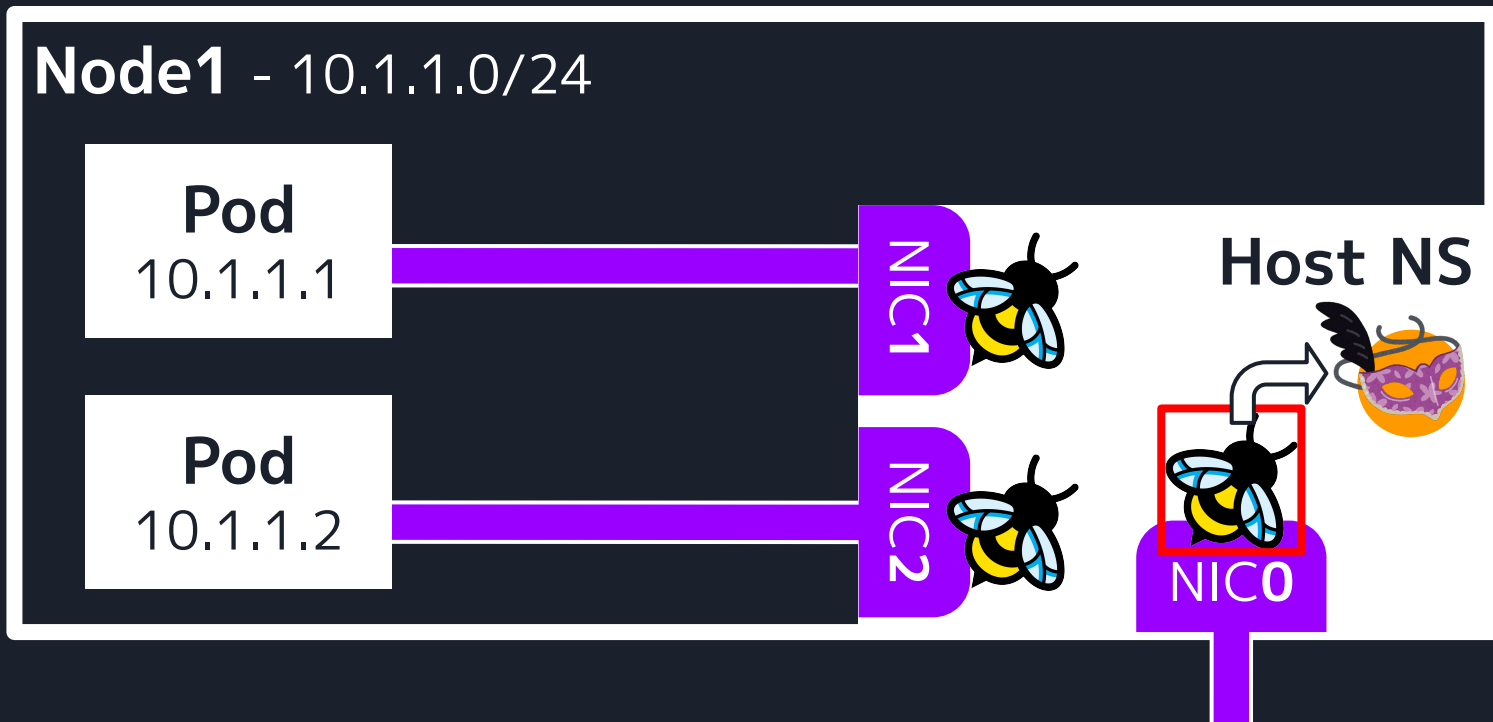
**Pod**  
10.1.1.2

NIC1

NIC2

**Host NS**

NIC0





# Cilium - 特徴まとめ

- **同じノード内でのL3疎通**
  - **eBPF**を使って高速にルーティングを行う
- **異なるノード間のL3疎通**
  - **Flannel / Calicoと同じ仕組みを使う**
  - 外からノードに入ってきたパケットは「同じノード内での疎通」と同じく**eBPF**でPodまで運ばれる
- **eBPFでIPマスカレードする**



## ところで: 汎用的 ≠ それだけでいい

- これらのOSSは**汎用**を目指している
  - どんな環境でも動くように、**一般的な仕組み**のみを利用して実装をしている
- メジャーなOSSプラグインは、どんな所でも使える**汎用性**の代わりに**オーバーヘッド**を許容している
  - 特に**オーバーレイ**は顕著に遅い
  - **SDN**はそれなりにノードに処理を要求する



# 環境によって効率の良い形は違う

- なぜCNIはKubernetesにデフォルトで組み込まれずわざわざ**インターフェース**になっているのか
  - クラウドベンダーが自社のネットワークに合わせて**効率の良い処理基盤**を組むため
- CNIの魅力は、実装が完全に**ユーザー依存**であること
  - オンプレ環境でネットワークが全部わかるなら、それに**効率化**するに越したことはない



# 環境特化なCNIのインセンティブ

- 既存のアンダーレイネットワークに最適化したい
  - 既にBGP/EVPN/VRF等、成熟した基盤があるとき
- 既存の運用・監視・障害対応フローに合わせたい
  - 既存の監視基盤に寄せられる形でメトリクス・イベント・トレースを最初から設計したい
- Podをもっとシームレスに繋ぎたい
  - 既存のネットワークと深く統合するとVMとPodがシームレスに繋がったりする





# 環境特化なCNIプラグインの例

- **GKE Dataplane V2**
  - **GKEでデフォルト**になっているプラグイン
- **AWS VPC CNI**
  - EKSで使える、**VPCをベース**にしたプラグイン
- **Coil on Neko**
  - **Cybozu**が開発しているプラグイン

それぞれの概要を軽く見ていきましょう

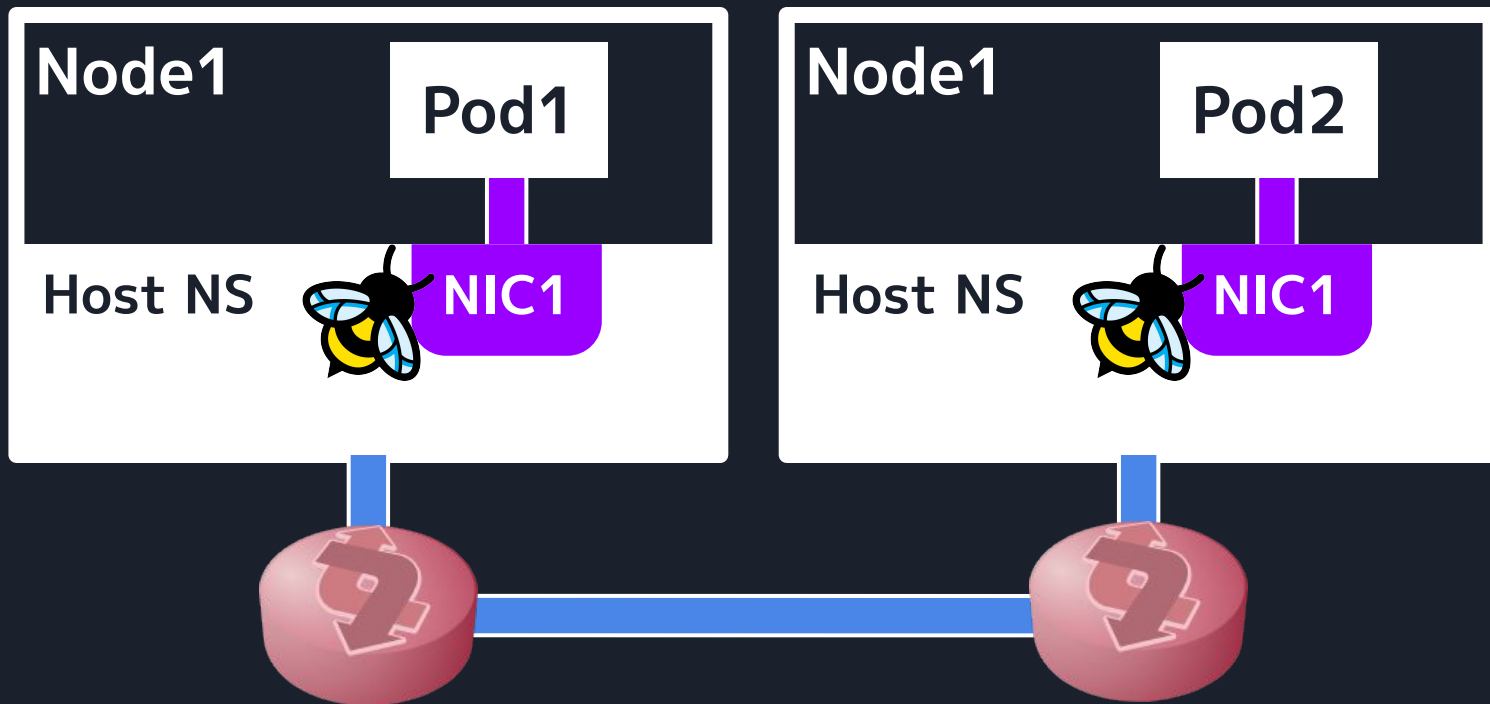


# GKE Dataplane V2

- GKEの開発チームが**Cilium**を**カスタマイズ**したもの
- **同じノード内**でのL3疎通
  - **Cilium**の仕組みをそのまま用いる
- **異なるノード間**のL3疎通
  - 準備: **アンダーレイネットワークに設定を流し込む**
  - 同ノード内のPod宛でなければ**そのまま外に出す**
  - **アンダーレイネットワークがルーティングする**

# GKE Dataplane V2 - 異なるノード間

Pod作成時: アンダーレイネットワークに設定を流し込む



# GKE Dataplane V2 - 異なるノード間

Pod作成時: アンダーレイネットワークに設定を流し込む

Node1

Pod1

Pod1、Pod2の  
経路情報をCNIが教える

Node1

Pod2

Host NS



NIC1



## GKE Dataplane V2 - 異なるノード間

このような仕組みで  
ただパケットを**外に出すだけで**  
勝手に**相手がいるノードに**  
**届く状態**が実現される

総



# Amazon VPC CNI

- **AWS VPCとの連携**に非常に強いCNI
- **VPCの中から**PodのIPアドレスが割り当てられる
  - ノードの固定サブネットが無いので、**アドレス空間の効率良い利用**ができる
- 疎通性の確保も**VPCを通して**行われる
  - 詳細はあまり公開されていない
- 同じVPC内にいれば、**VM等と直接通信**ができる

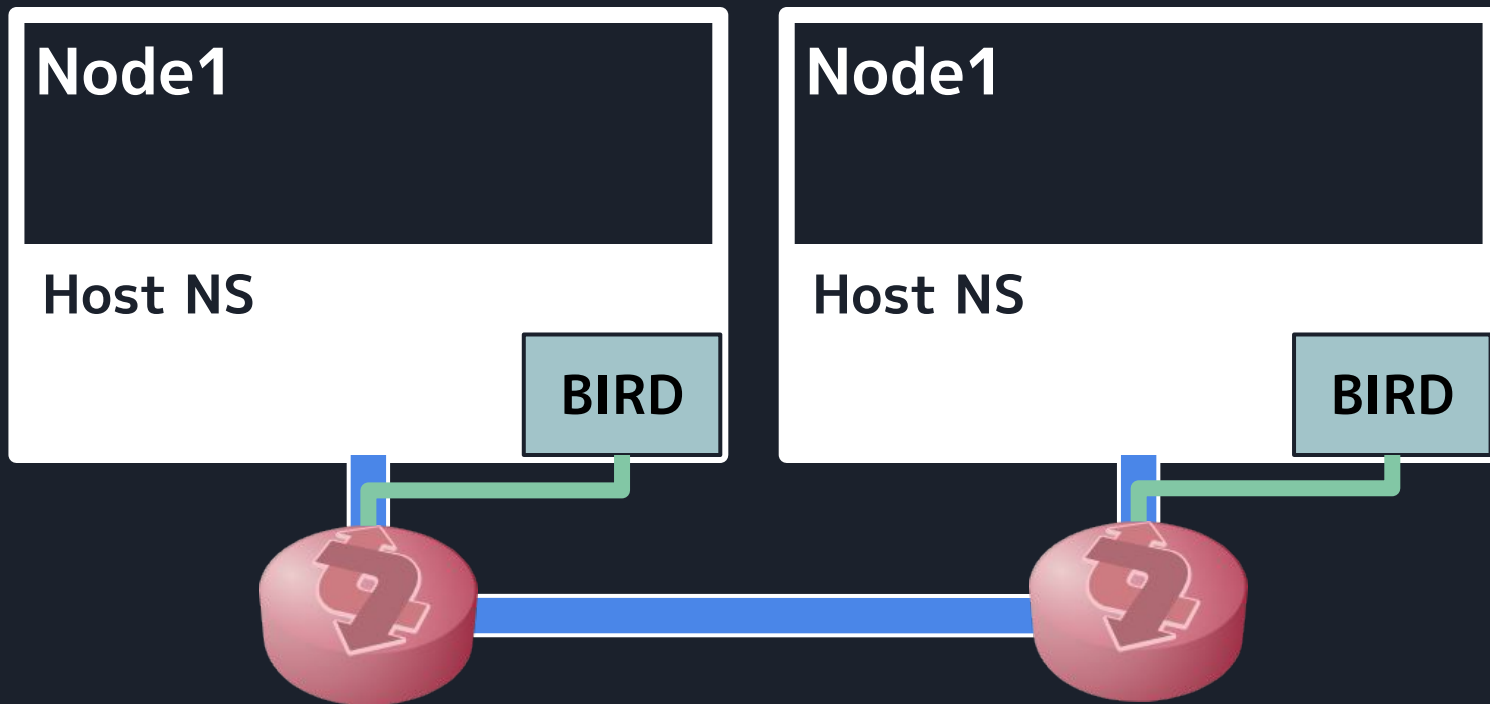


# Coil

- Cybozuで**Neko基盤**を作っているチームが開発している  
**オープンソース**のCNI
- Neko基盤内では、**Cilium** & **BIRD**と組み合わせて運用
- **同じノード内**でのL3疎通
  - **Cilium**の仕組みをそのまま用いる
- **異なるノード間**のL3疎通
  - Calicoと同じくBIRDから**BGPで経路広報**するが、  
**アンダーレイネットワーク**に経路を流す

# Coil - 異なるノード間

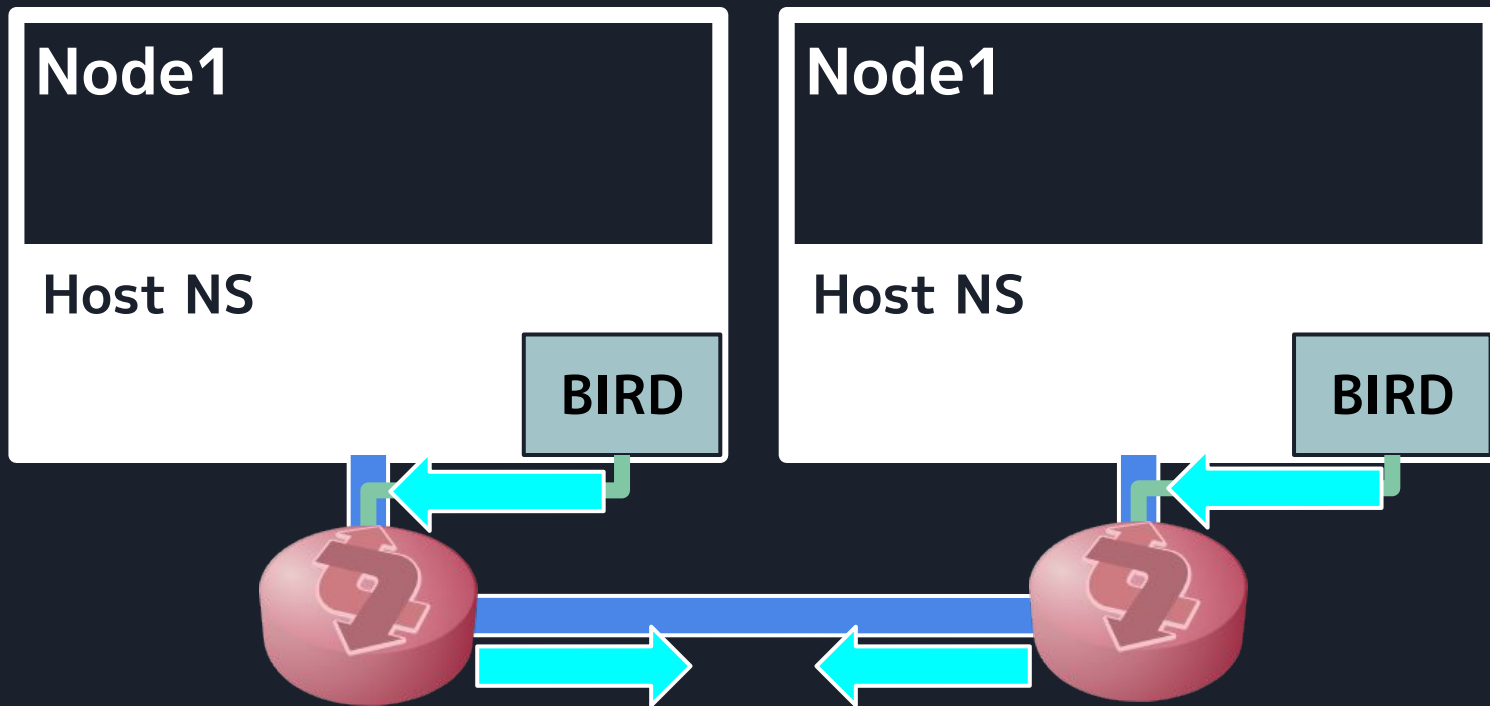
下準備: BIRDでアンダーレイネットワークとBGPピアを張る






# Coil - 異なるノード間

下準備: Podサブネットを全体で互いに**広報**する





CNIはどんな実装も受け入れる  
**柔軟な思想**の表れであり  
**様々な実装**がある

皆さんも馴染みのある技術が多かったのでは



# アジェンダ

- イントロダクション
- Kubernetesを知る
- Kubernetesとネットワーク
- Linuxのネットワークスタック
- L2/L3ネットワーク: CNI ←イマココ
- L4ロードバランサー: kube-proxy
- プラットフォームとの対話で広がる可能性
- まとめ



# アジェンダ

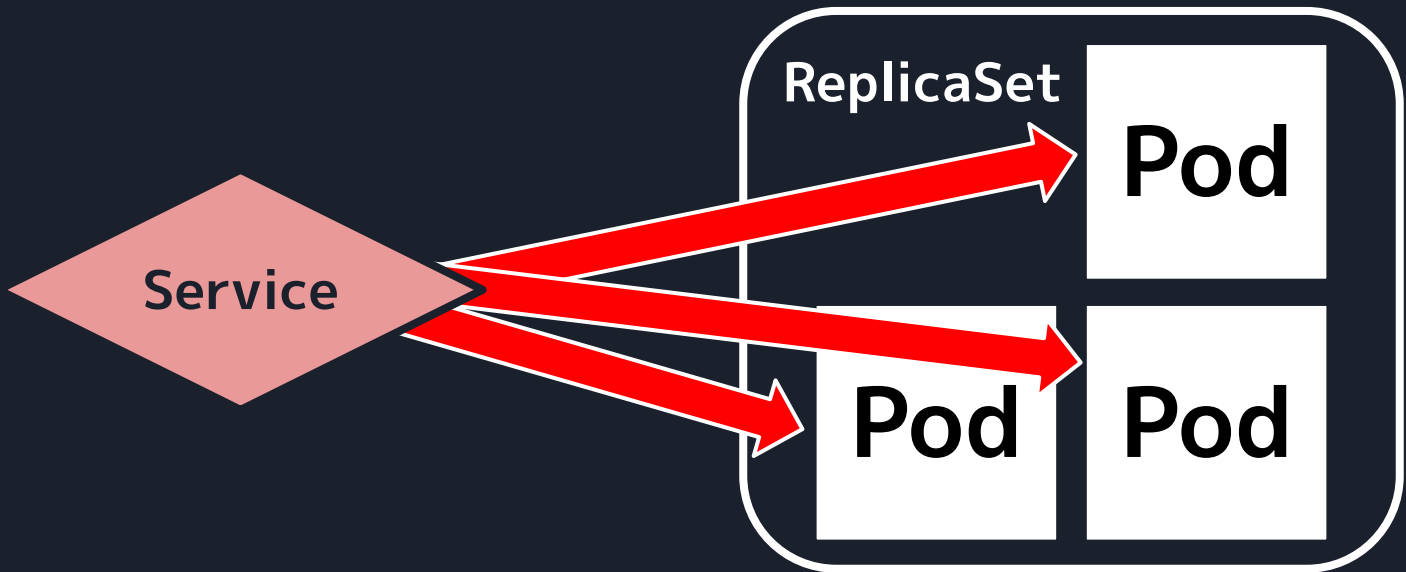
- イントロダクション
- Kubernetesを知る
- Kubernetesとネットワーク
- Linuxのネットワークスタック
- L2/L3ネットワーク: CNI
- L4ロードバランサー: kube-proxy ←イマココ
- プラットフォームとの対話で広がる可能性
- まとめ

# L4ロードバランサー: **kube-proxy**



## 【復習】 Service

- 同じ内容のPod群に、均等に**通信を振り分ける** (L4LB)
- Pod群を**1つのIPアドレス**でまとめることができる





# kube-proxy

- **Service**を担当するコンポーネント
- Kubernetesが**デフォルト**で持っている
- kube-proxyのやる事
  - **宛先**IPアドレスが**Service**の通信があったら、所属する任意PodのIPで**DNAT**する
  - 書き換え先のPodの**バランス**を**いい感じ**にする



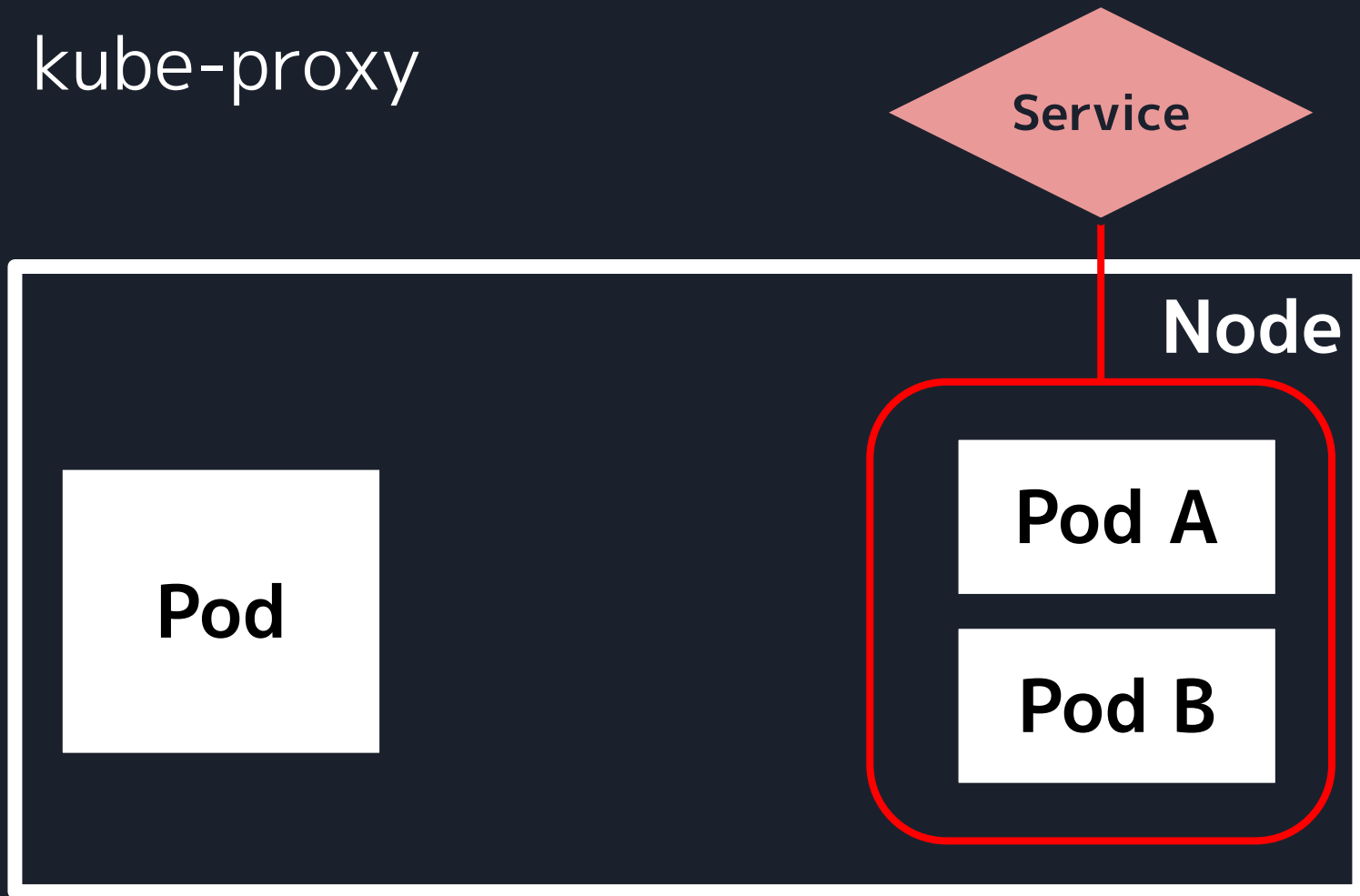
# kube-proxyの実装方式

- 基本的には**netfilter**を利用してNATしている
  - kube-proxyは保存されているServiceの定義を監視して**NATルールを流し込む**だけ
  - 実際にプロキシしているわけではない
- 現在のデフォルト: **iptables**を通して設定
- これから: **nftables**への移行が進む (現在ベータ)



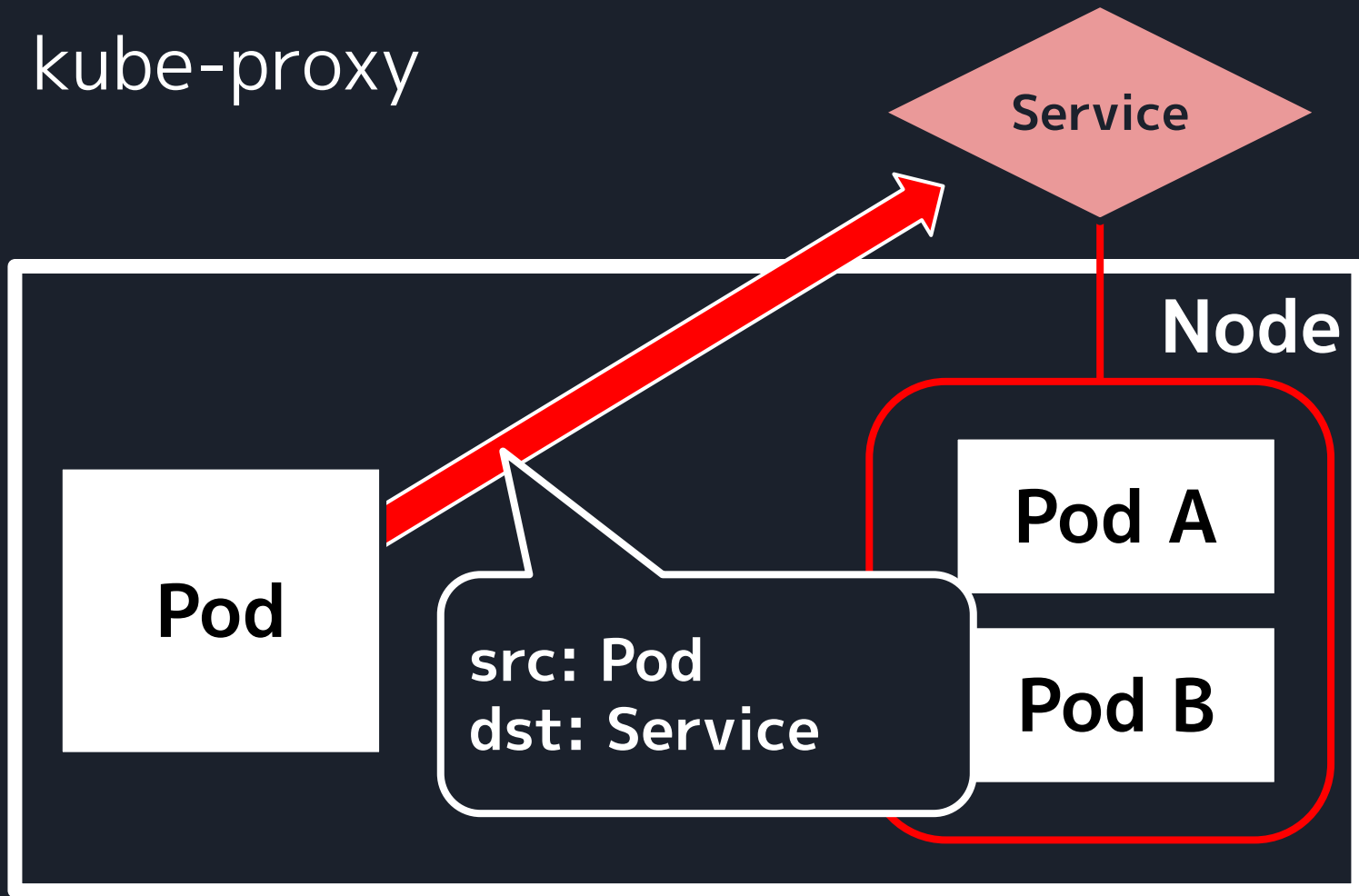


# kube-proxy

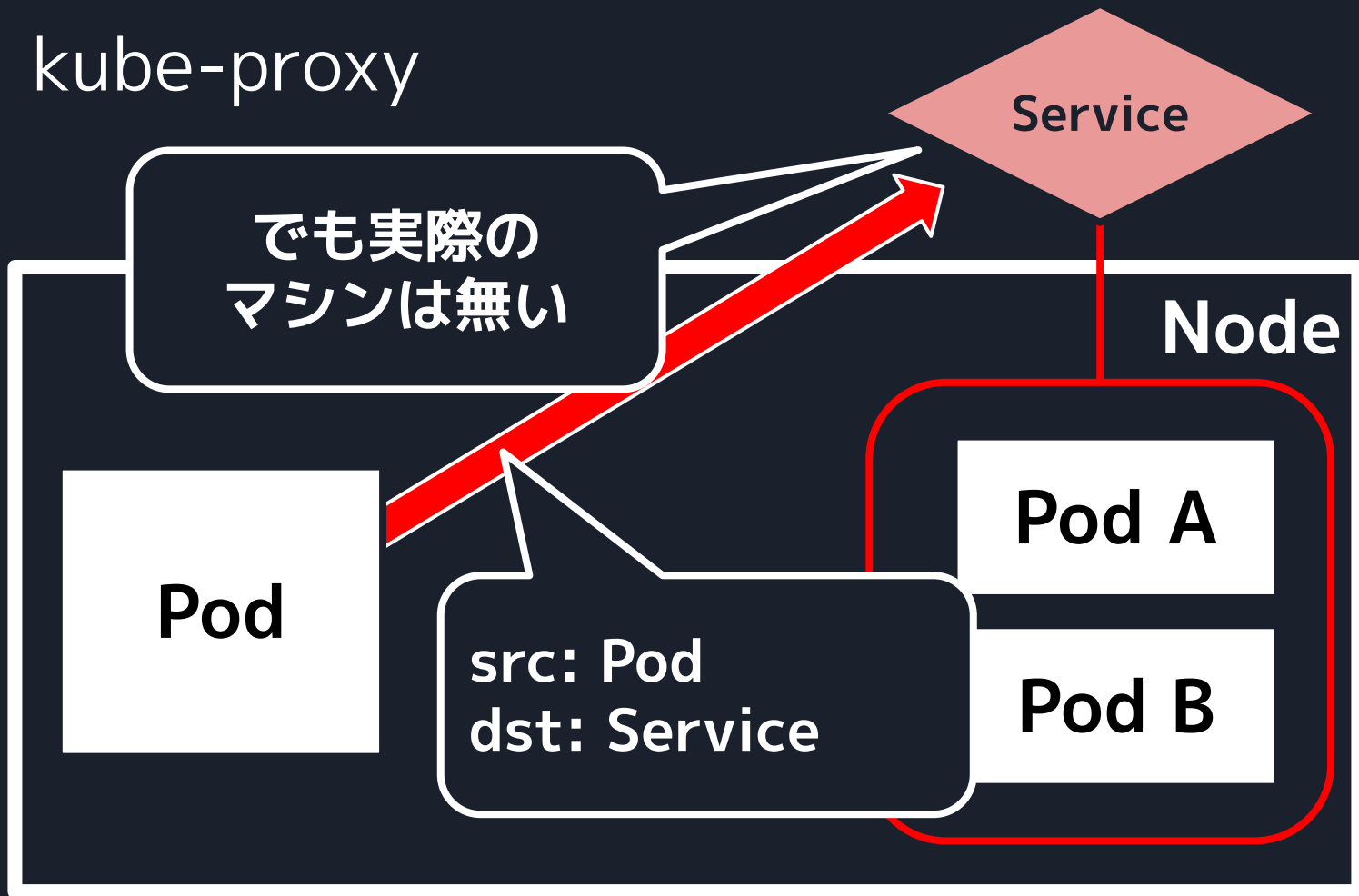




# kube-proxy



kube-proxy



kube-proxy

なんとか  
します！

Service

Node

netfilter

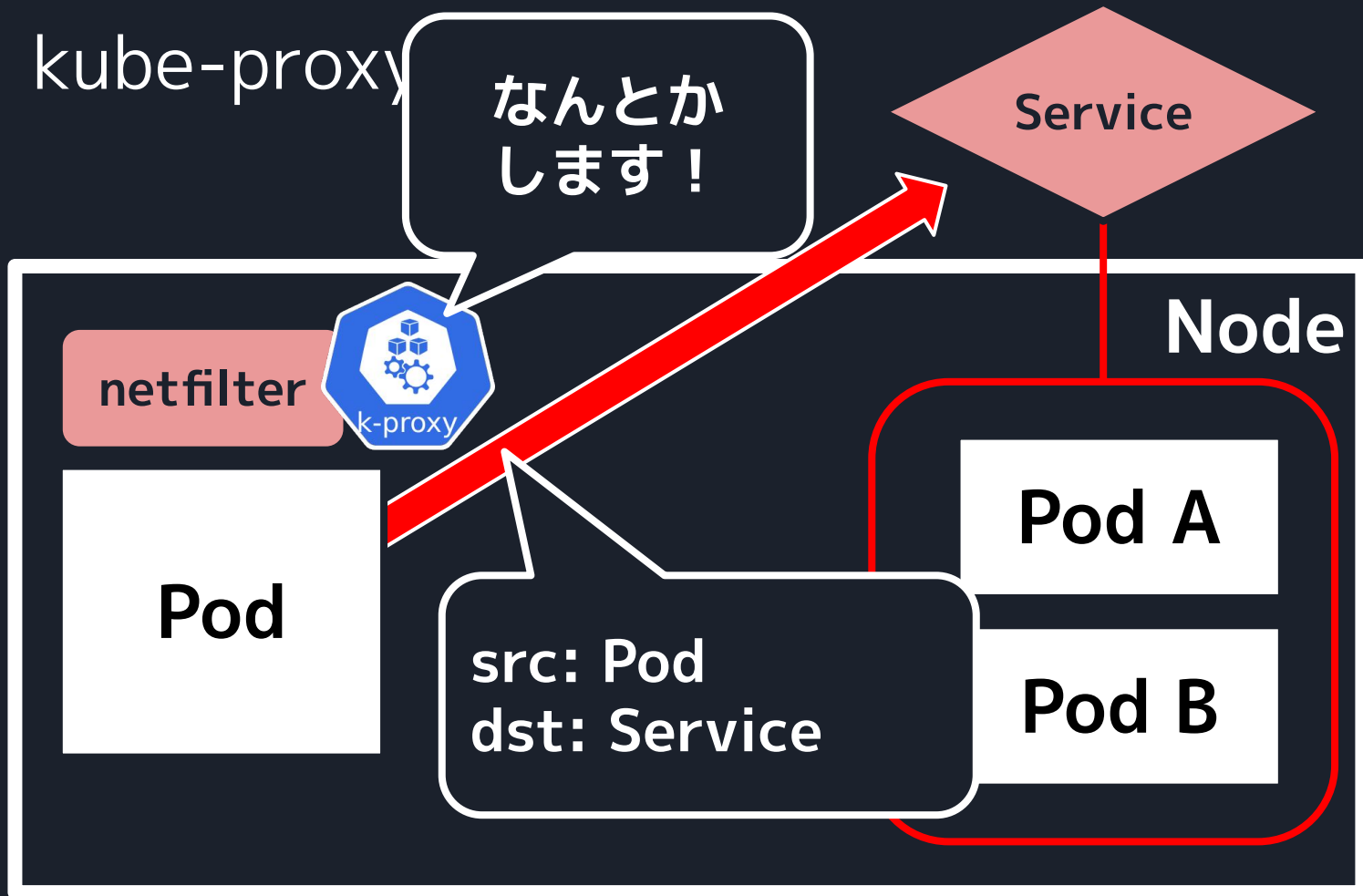
k-proxy

Pod

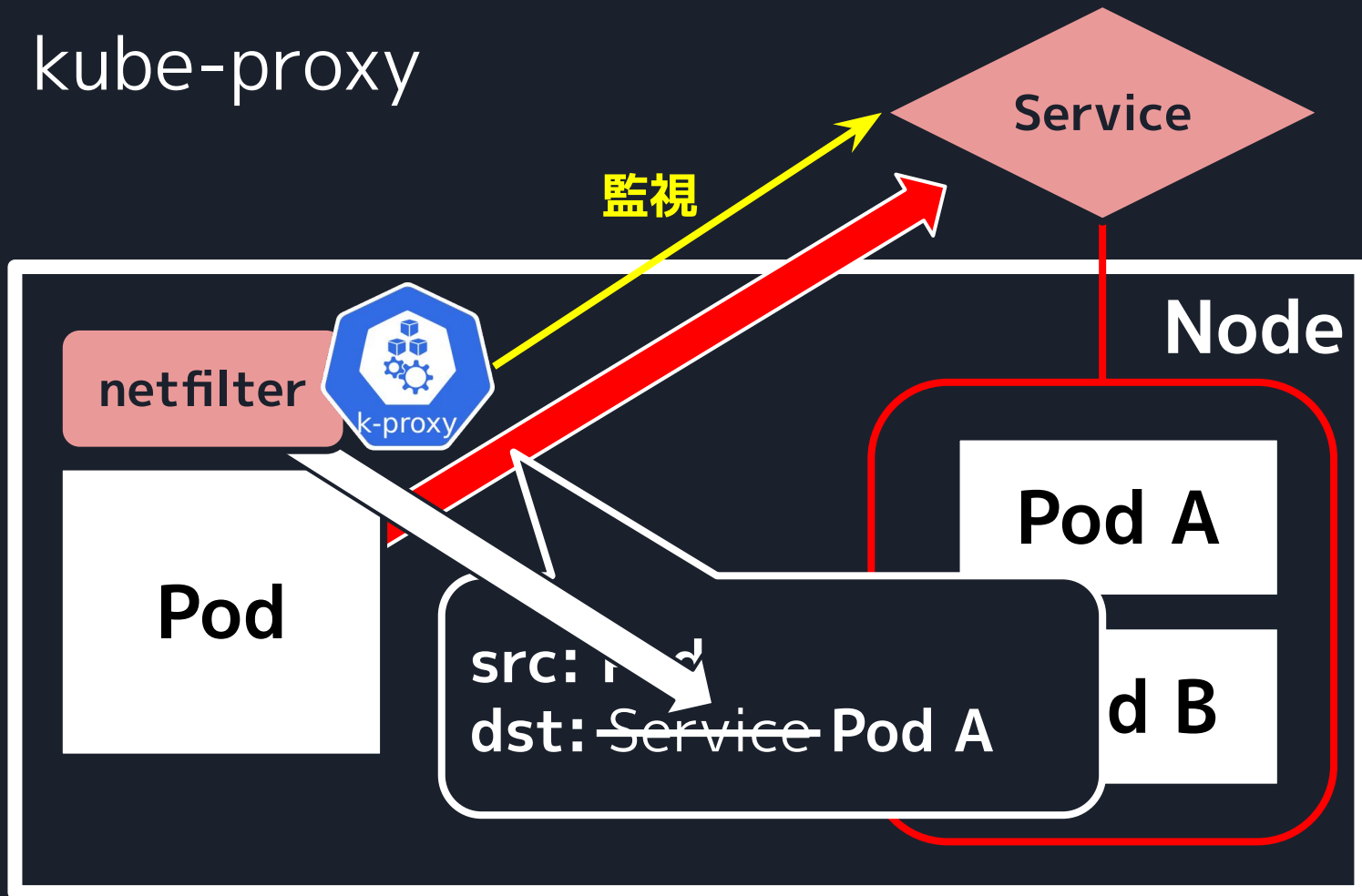
Pod A

Pod B

src: Pod  
dst: Service

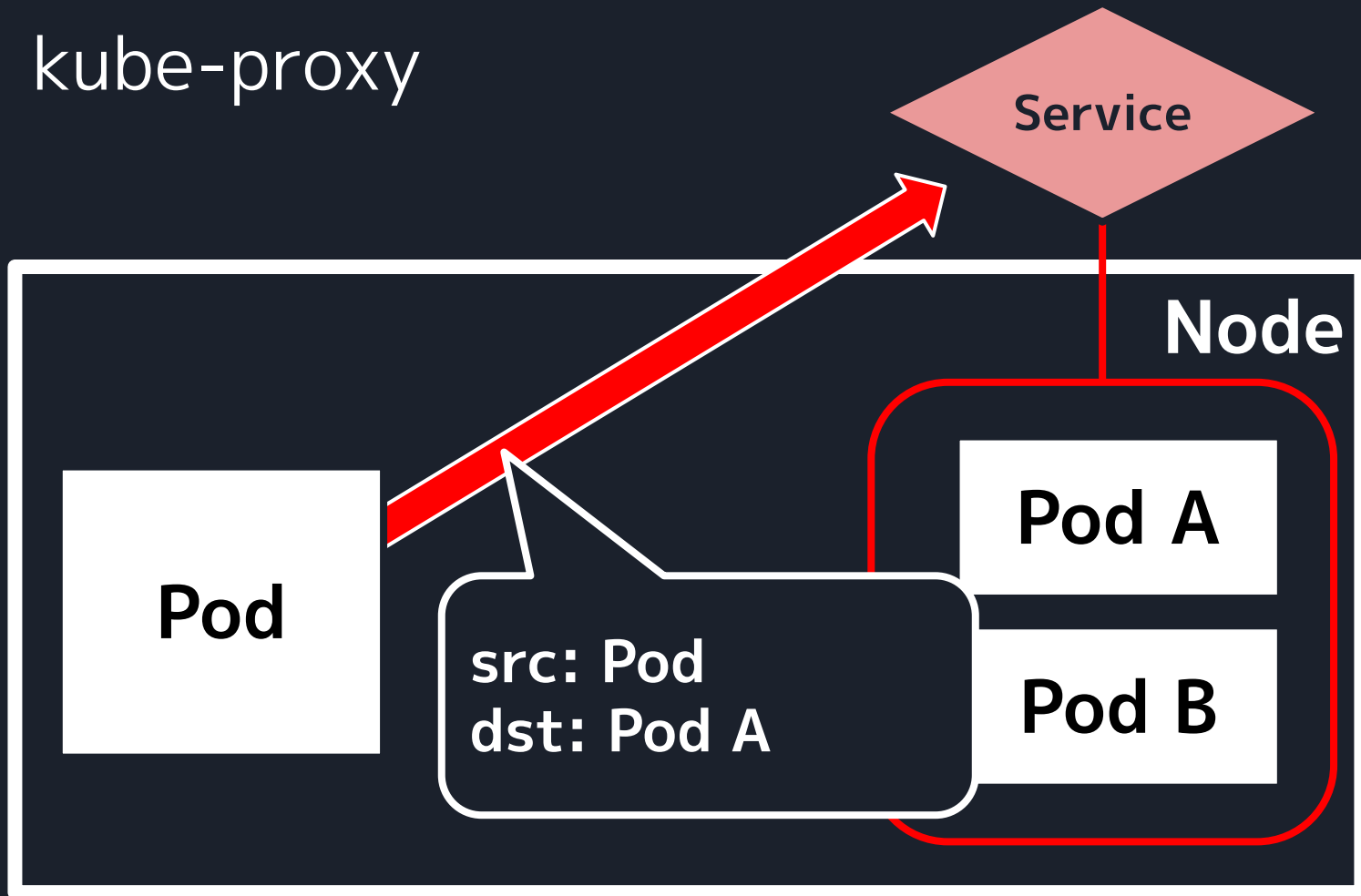


kube-proxy

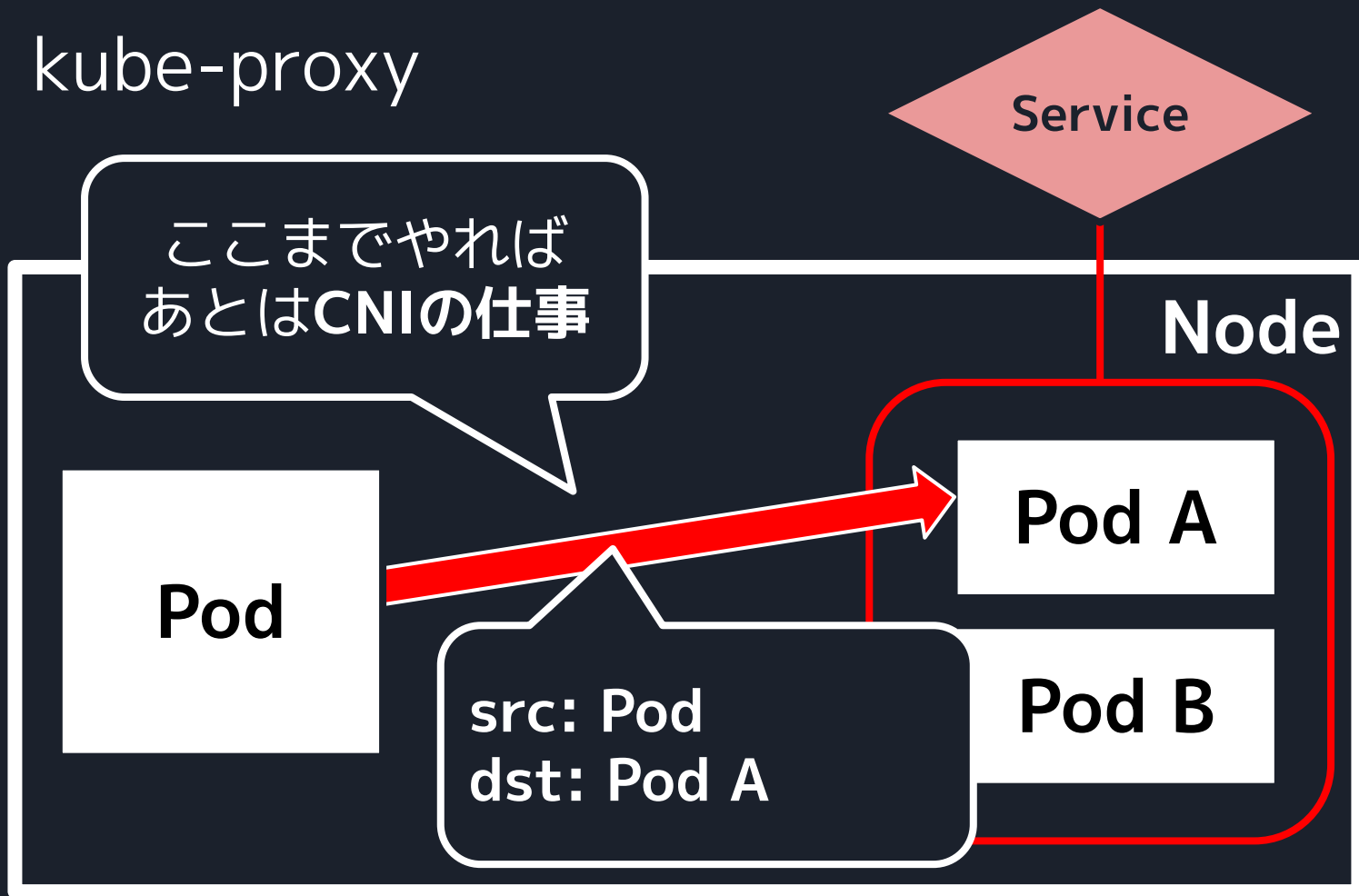




# kube-proxy

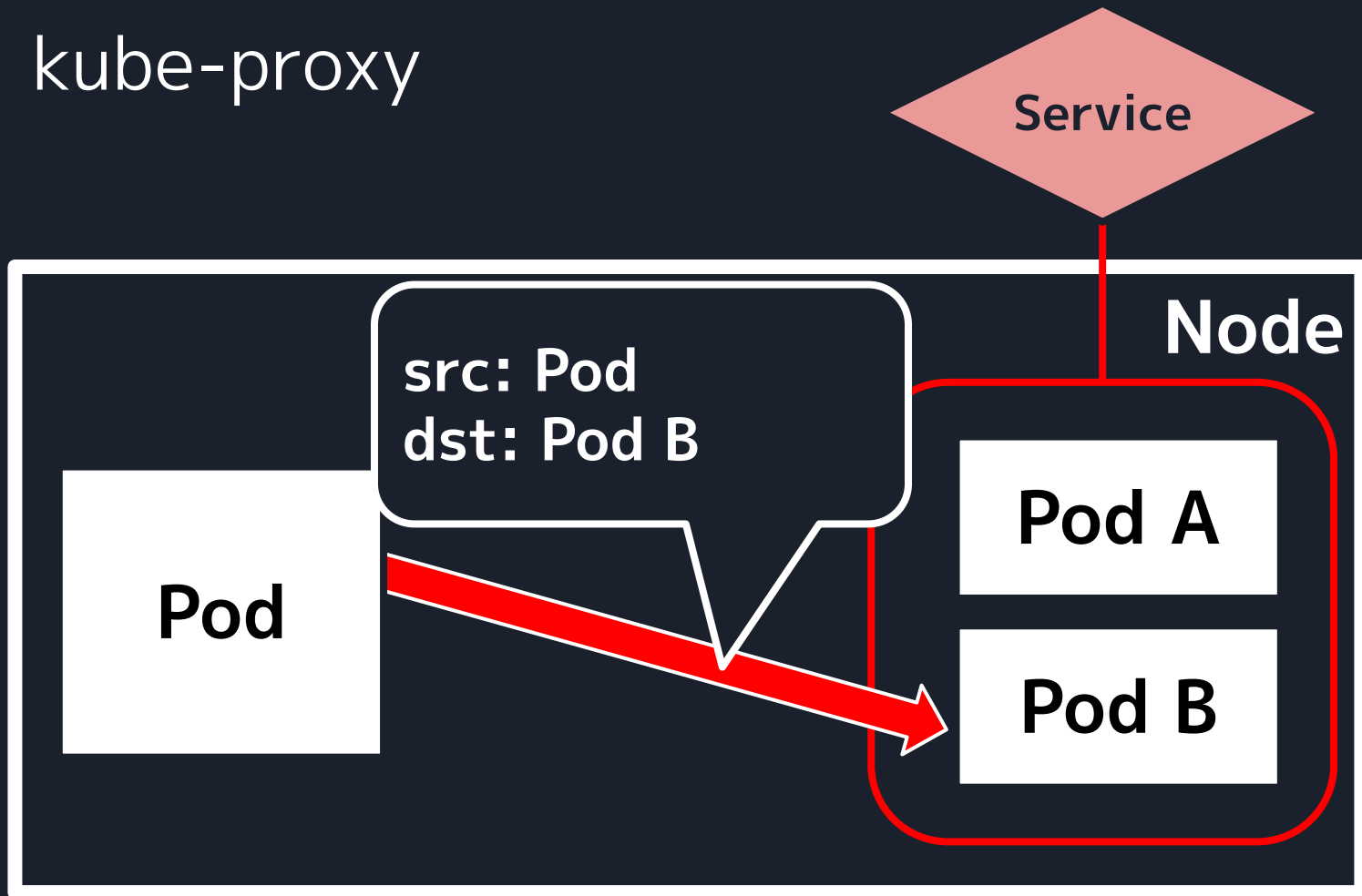


kube-proxy





# kube-proxy







kube-proxy

Service

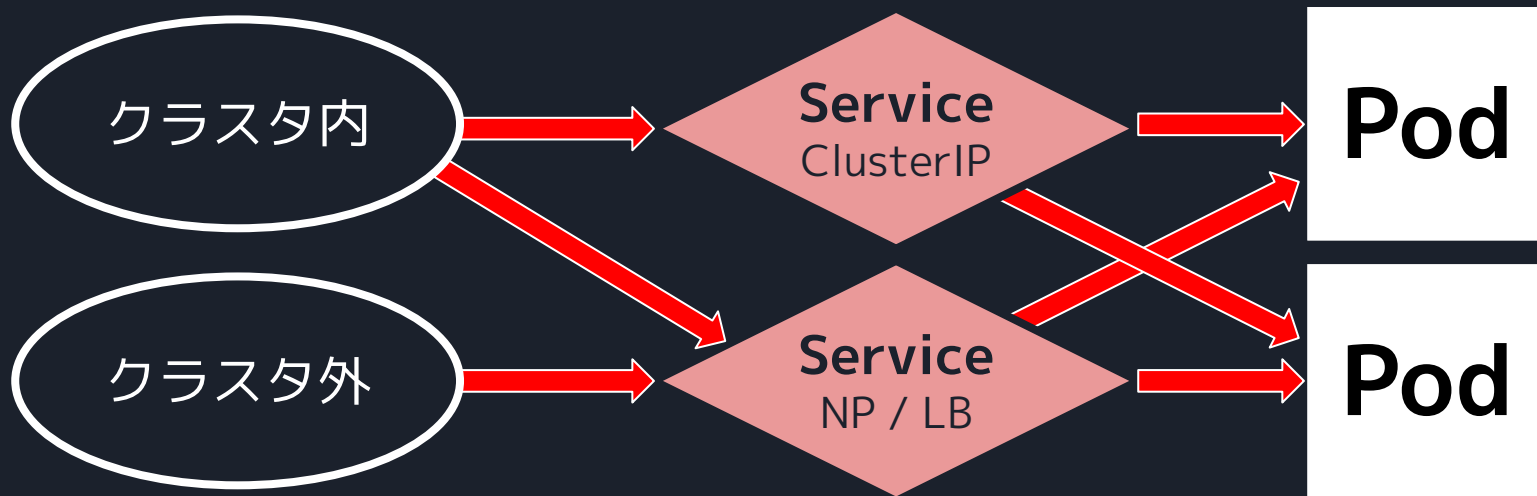
kube-proxyの役目は  
Serviceに属するどこかの  
Podに**DNAT**するだけ



Pod B

## 【復習】 Serviceのタイプ

- ClusterIP: **クラスタ内**からの通信のみ受け付ける
- NodePort / LoadBalancer: **クラスタ外**からの通信も受け付ける





# 外部から通信を受け付けるService

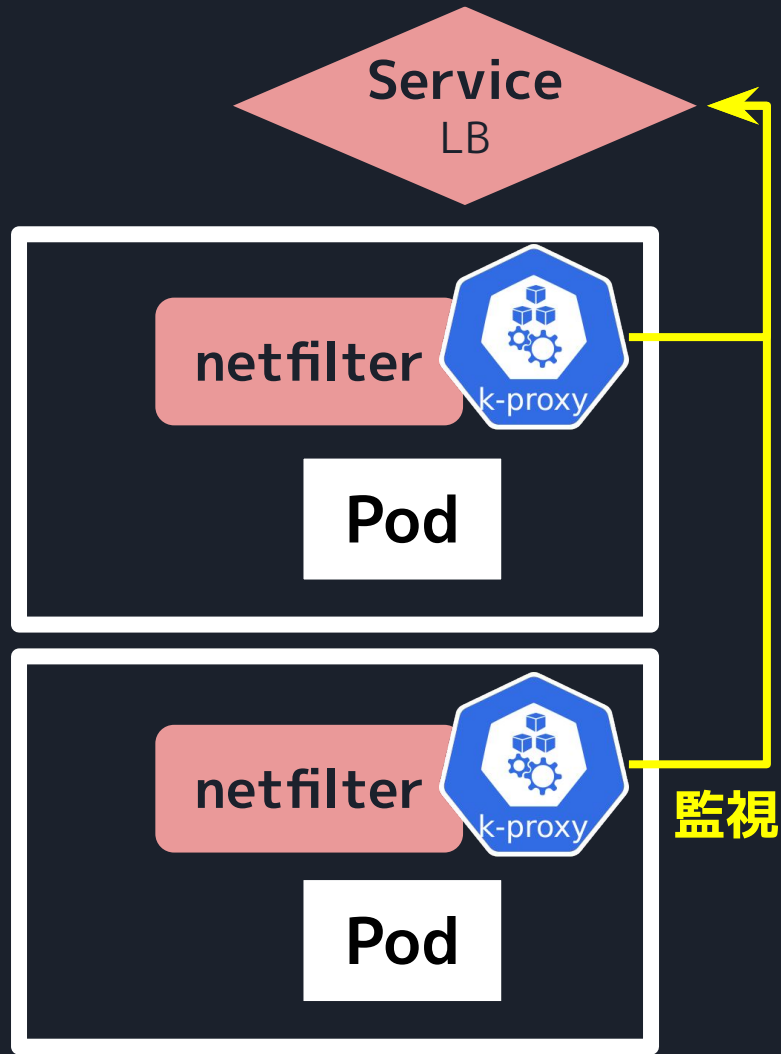
- **NodePort**

- 30000-32767のうち**1つのポート**を割り当て、**全ノードでそのポート宛**の通信を受け取る

- **LoadBalancer**

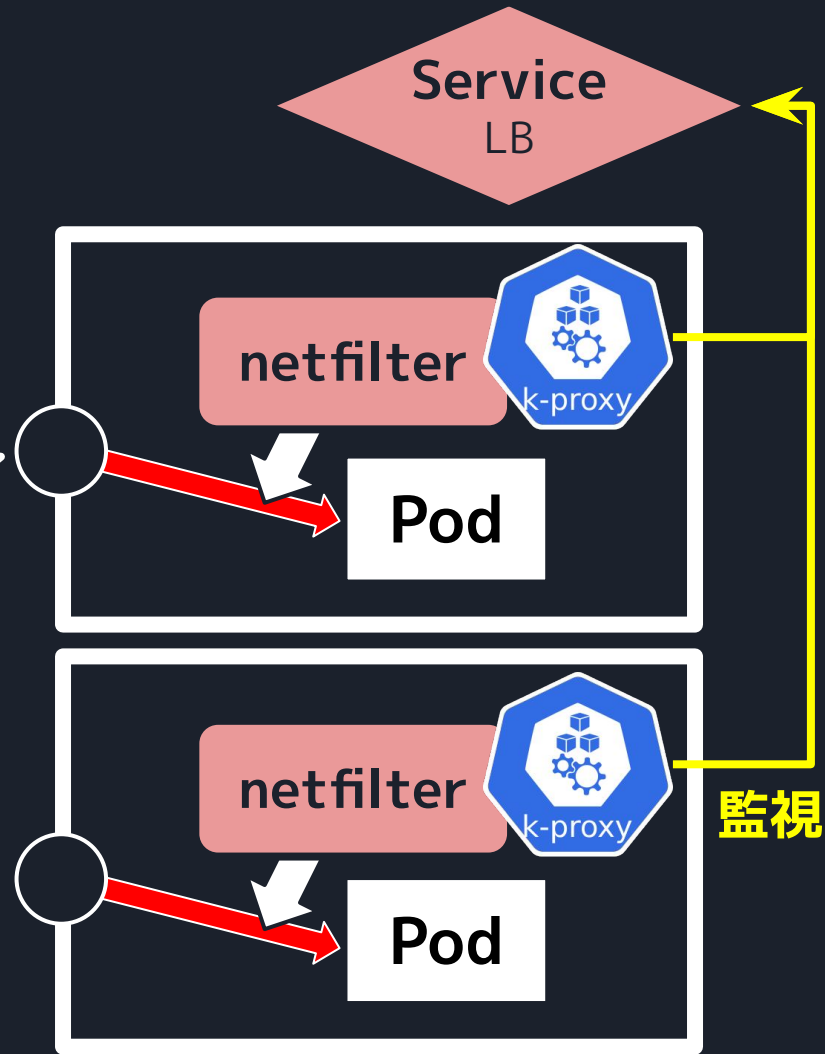
- Nodeportと同じ環境をそろえる
- **外部L4LBと連携しクラスタ外からアクセスできるIPアドレス**を割り当て、LBから ↑ で用意したポートにロードバランシングする

Type: NodePort



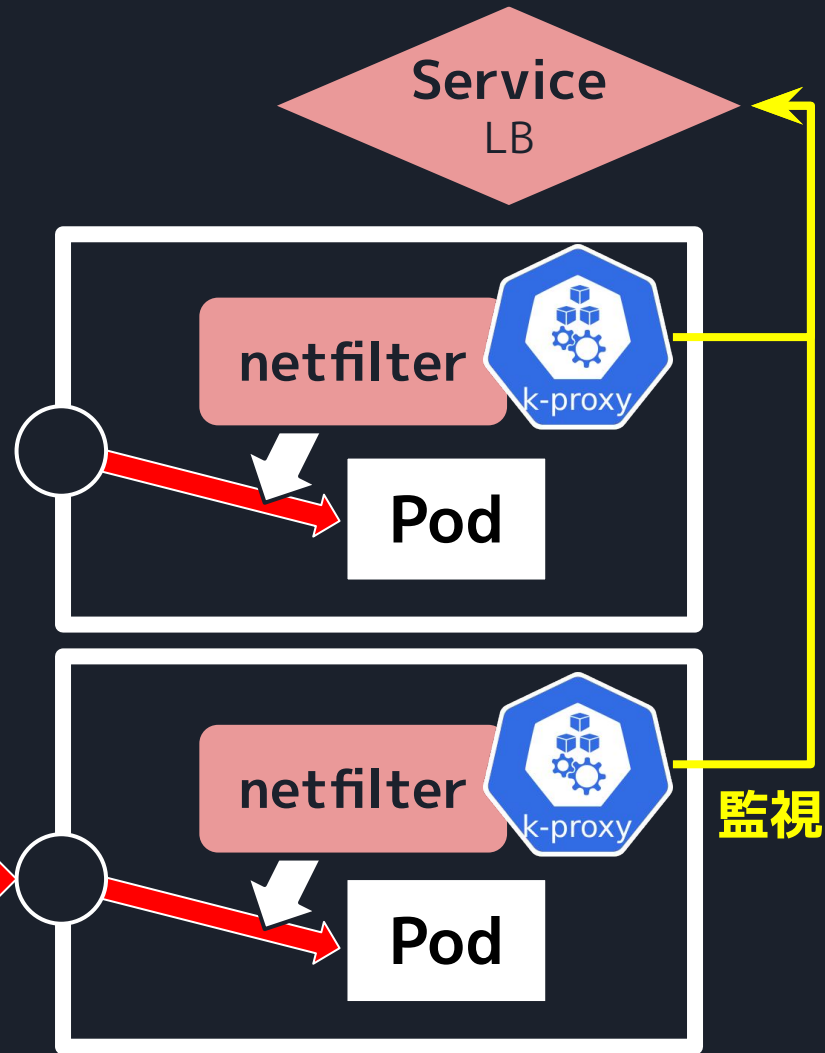
# Type: NodePort

ポートを1つ  
割り当て  
そこからの**通信**を  
**キャッチ**する



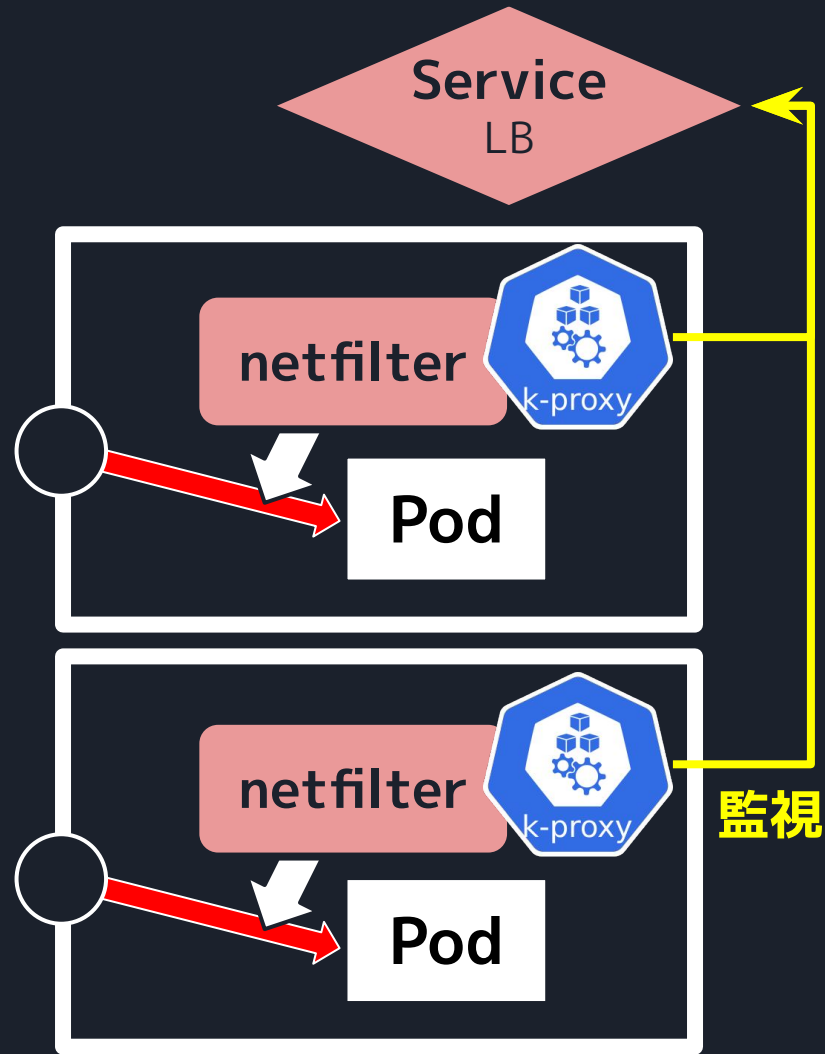
Type: NodePort

クラスタ外

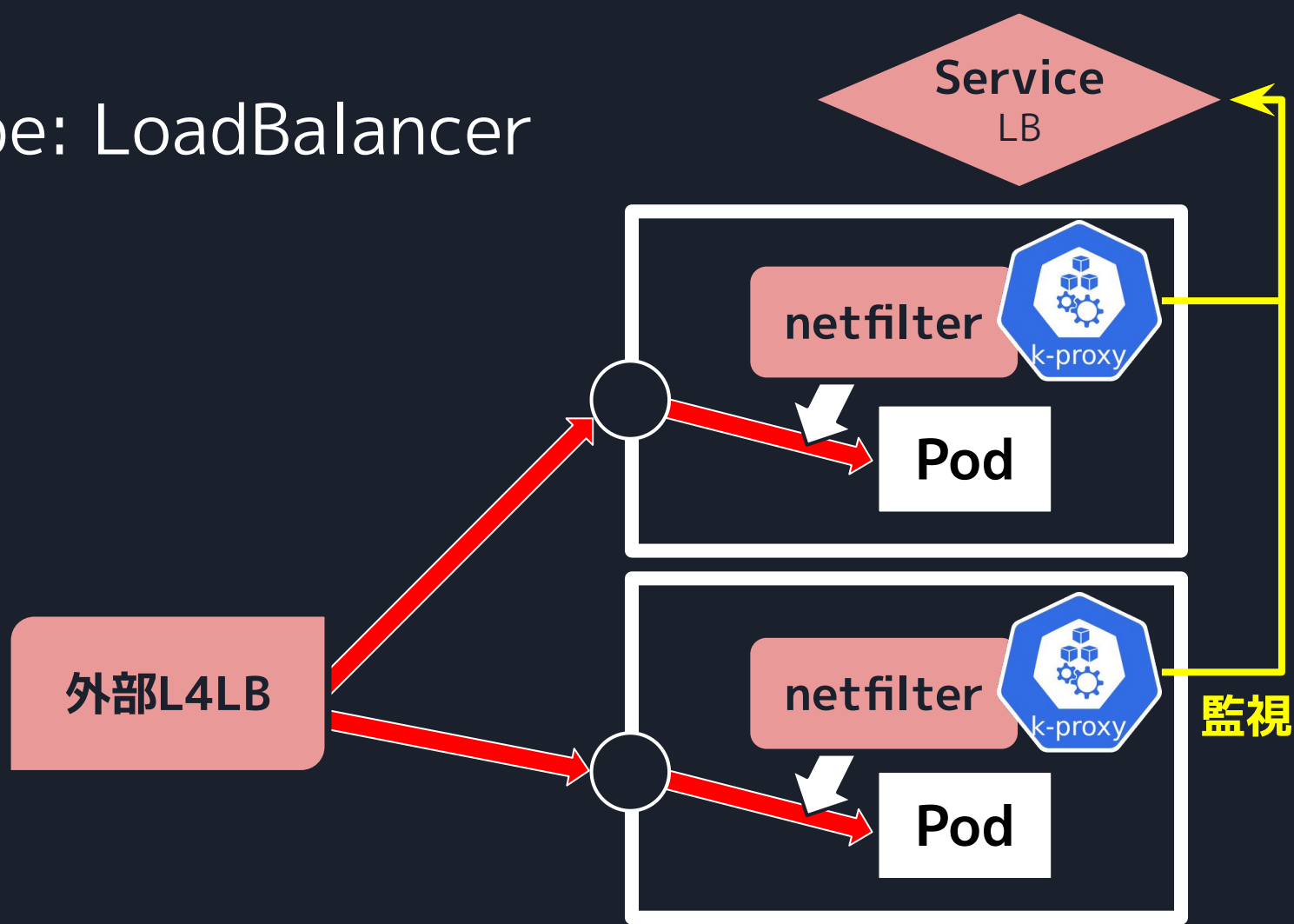


# Type: LoadBalancer

外部L4LB



# Type: LoadBalancer





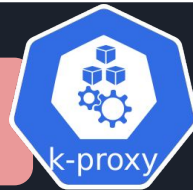
# Type: LoadBalancer

クラスタ外から  
アクセス可能な  
IPアドレスを  
割り当て

外部L4LB

Service  
LB

netfilter



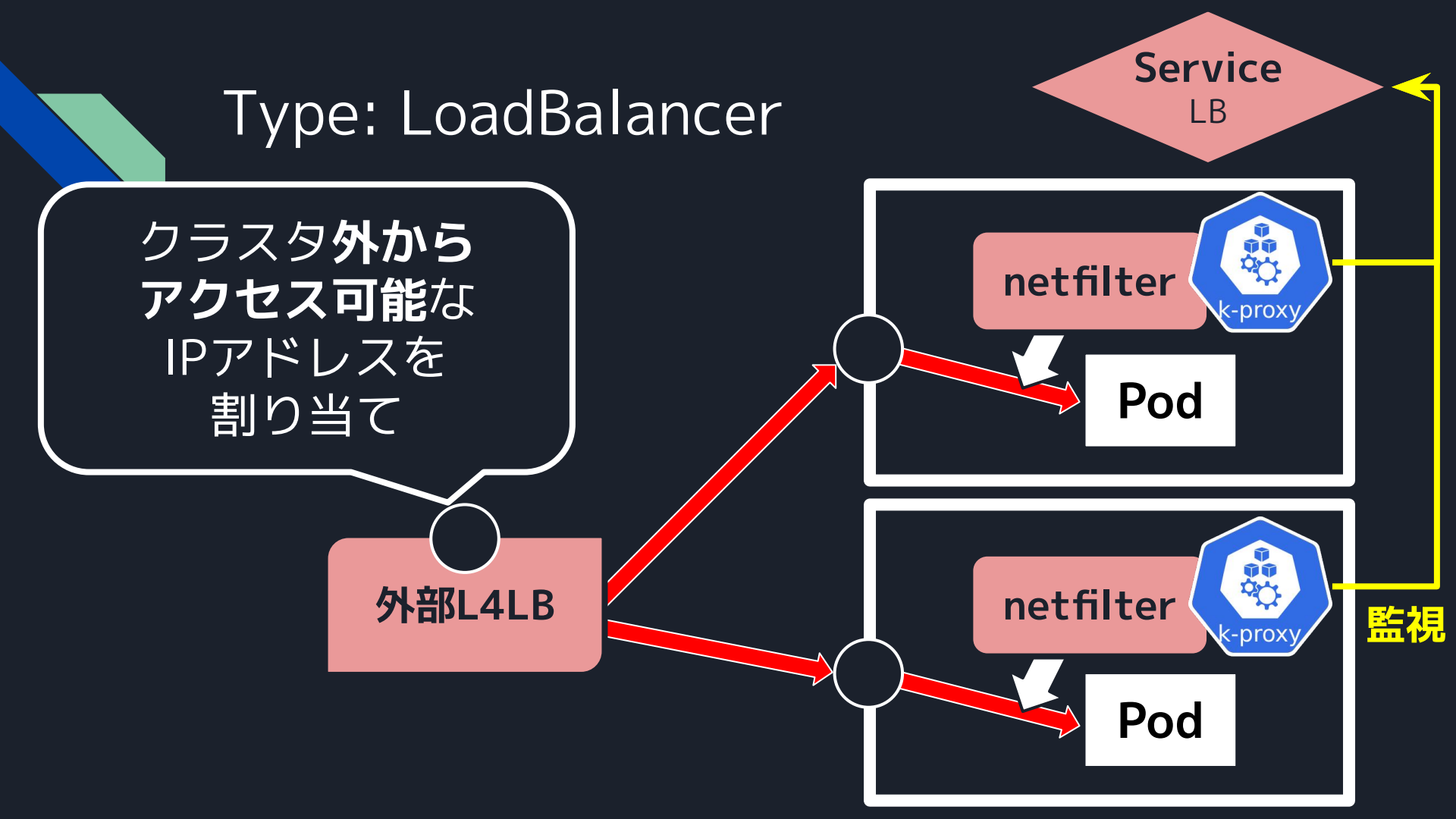
Pod

netfilter

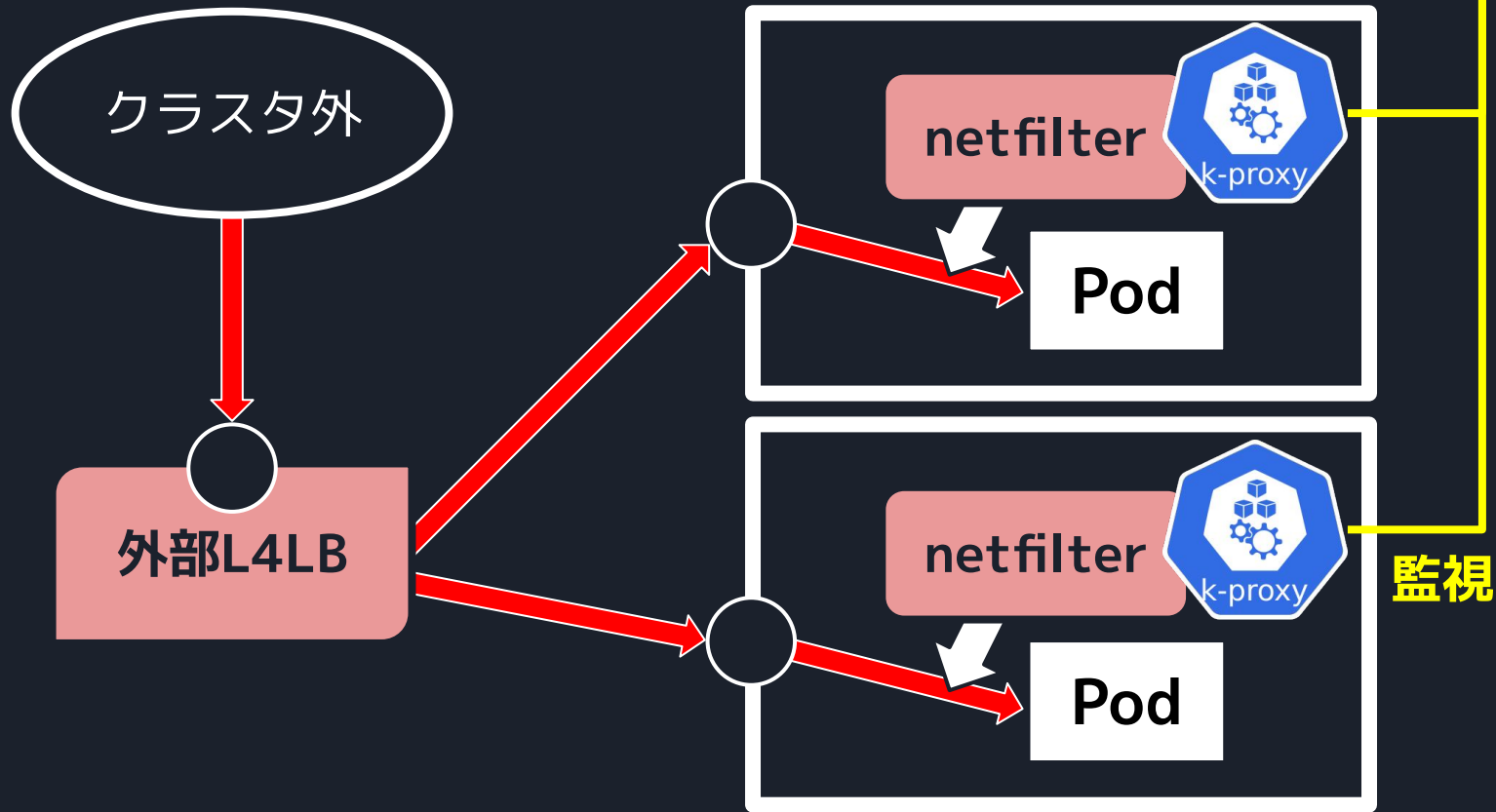


Pod

監視



# Type: LoadBalancer





# 外部L4LB

- Type: LoadBalancer用の**外部L4LB**は、Kubernetesがデフォルトで用意していないコンポーネント
  - 既存の基盤に合わせて**選定** or **自作**の必要がある
- 実装例
  - **MetalLB**
    - オンプレ環境を想定した**汎用OSS**
  - **AWS Load Balancer Controller**
    - AWSの**NLB**を使った実装



# アジェンダ

- イントロダクション
- Kubernetesを知る
- Kubernetesとネットワーク
- Linuxのネットワークスタック
- L2/L3ネットワーク: CNI
- L4ロードバランサー: kube-proxy ←イマココ
- プラットフォームとの対話で広がる可能性
- まとめ



# アジェンダ

- イン트로ダクション
- Kubernetesを知る
- Kubernetesとネットワーク
- Linuxのネットワークスタック
- L2/L3ネットワーク: CNI
- L4ロードバランサー: kube-proxy
- プラットフォームとの対話で広がる可能性 ←イマココ
- まとめ


# プラットフォームとの 対話で広がる可能性





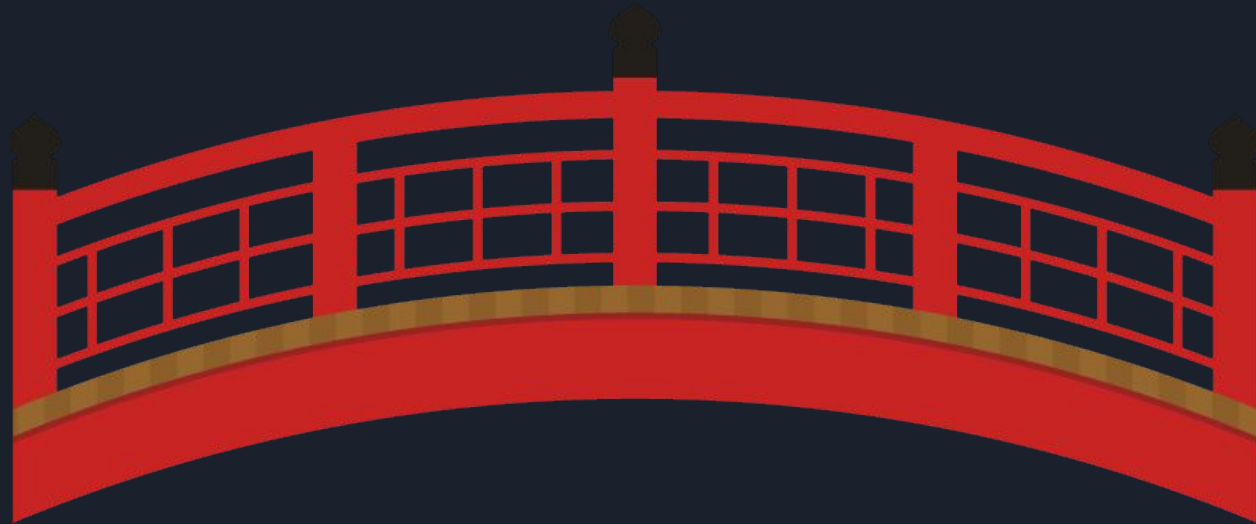
# 環境特化なネットワークを選択肢に！

- メジャーなOSSは、どんな所でも使える**汎用性**の代わりに**オーバーヘッド**や**複雑性**を許容している
  - **SDN**はそれなりにノードに処理を要求する
  - **設定項目**も多くなりがち
- K8sのネットワークは**差し替え可能**な部分が多い
  - **自分たちの**ネットワークの状態に合わせて、**最適なものを選ぶ/作**ることができる！
- この実現のためには、**対話が不可欠**



# プラットフォームと対話しよう

- 特にCNIはお互いの**責任を分割**するツール
- ブラックボックスから対話のための**共通言語**へ
- 「どちらも触れない」から「**どちらも触る**」にしたい





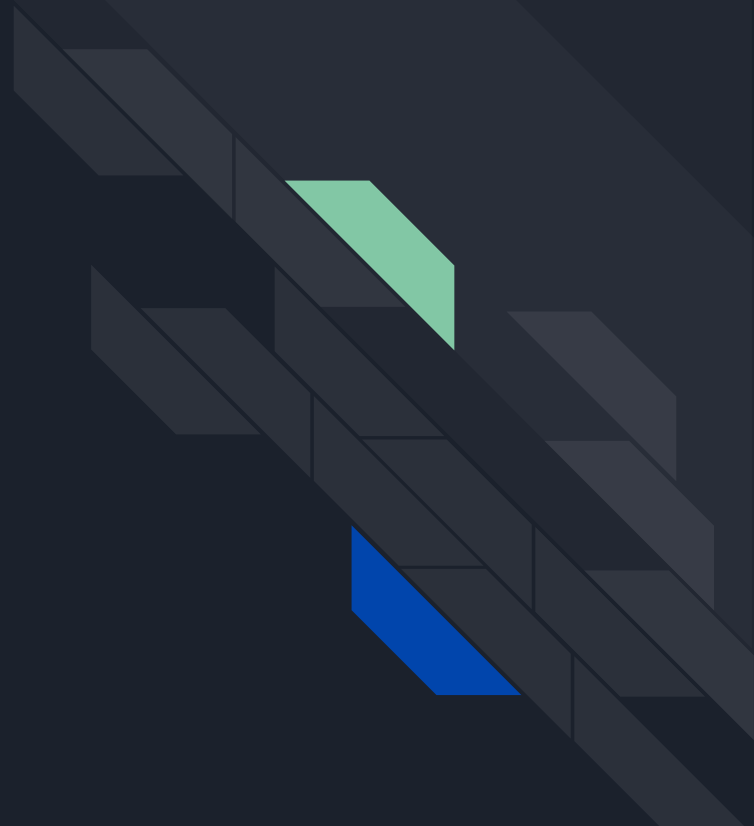
# 実はプラットフォーム方面の活動も…

- 昨年11月にあった、Kubernetes等クラウドネイティブなインフラに関するカンファレンスCNDW2025にて「**CNI徹底解説**」というタイトルでNagamiが登壇
- プラットフォームとネットワーク**双方から「共通言語」**を広めるのが目標
  - 何か良いアイデアがあればぜひ教えて下さい



**CLOUDNATIVE DAYS**  
**WINTER 2025**


まとめ





# 今日はこんなお話をしました

- イン트로ダクション
- Kubernetesを知る
- Kubernetesとネットワーク
- Linuxのネットワークスタック
- L2/L3ネットワーク: CNI
- L4ロードバランサー: kube-proxy
- プラットフォームとの対話で広がる可能性
- まとめ



# プラットフォームとの 「フラット」な対話で 共にスケールする基盤を！

もっとたくさんの詳細があります、入口になれば幸い



ありがとう  
ございました

参考になれば幸いです



## 【最後に】 議論ポイント

- 我々はプラットフォーム側の人間なので、**ネットワークエンジニアの視点**を分かり切れていない
- **ネットワークとプラットフォームはどのように関われるか**について色々な意見が聞きたい
  - ex) PF側・NW側視点でのCNI選定基準
  - ex) CNI改造・自製の判断基準、協力の仕方
  - ex) 責任分界点はどこに置くべきか
  - ex) コンテナネットワークの未来のビジョン