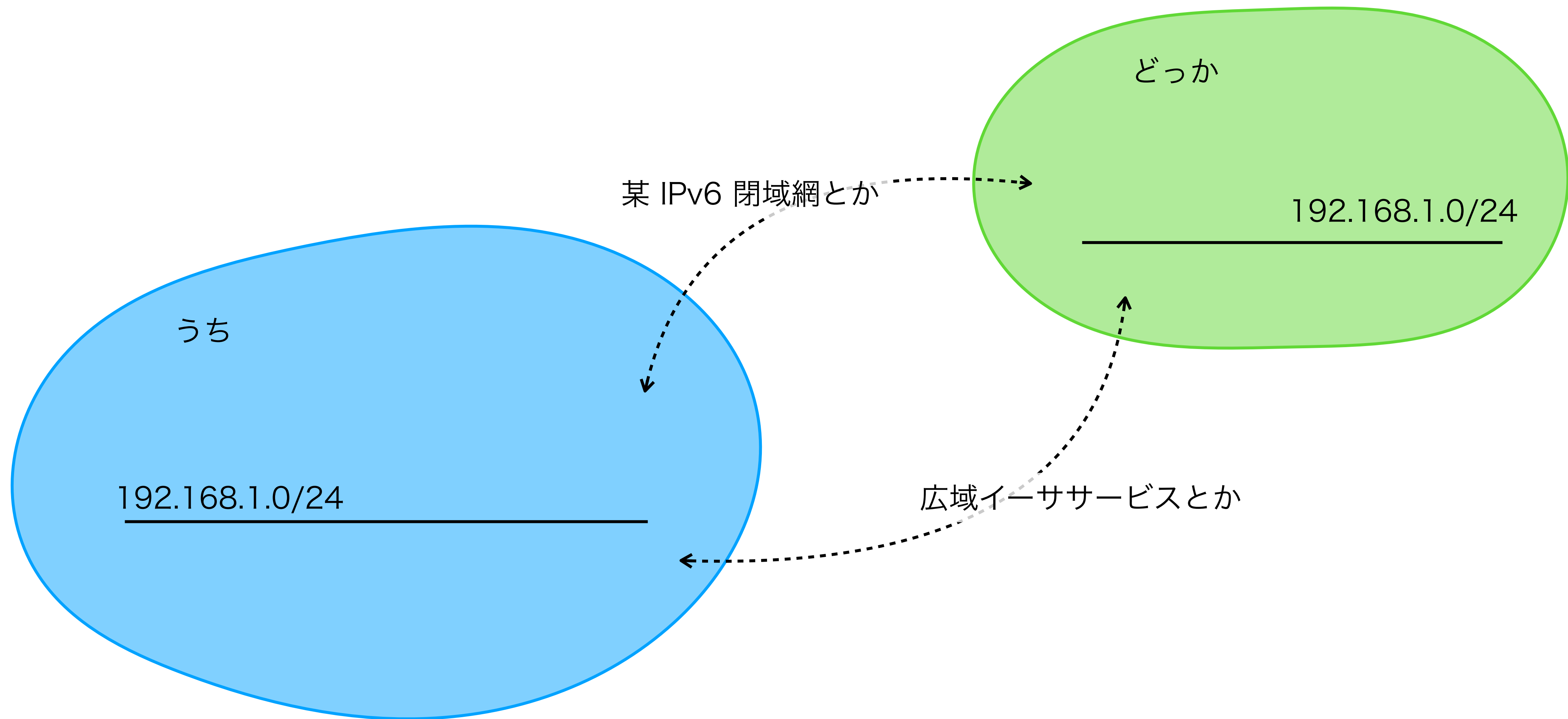


貧乏人のための L2 pseudowire

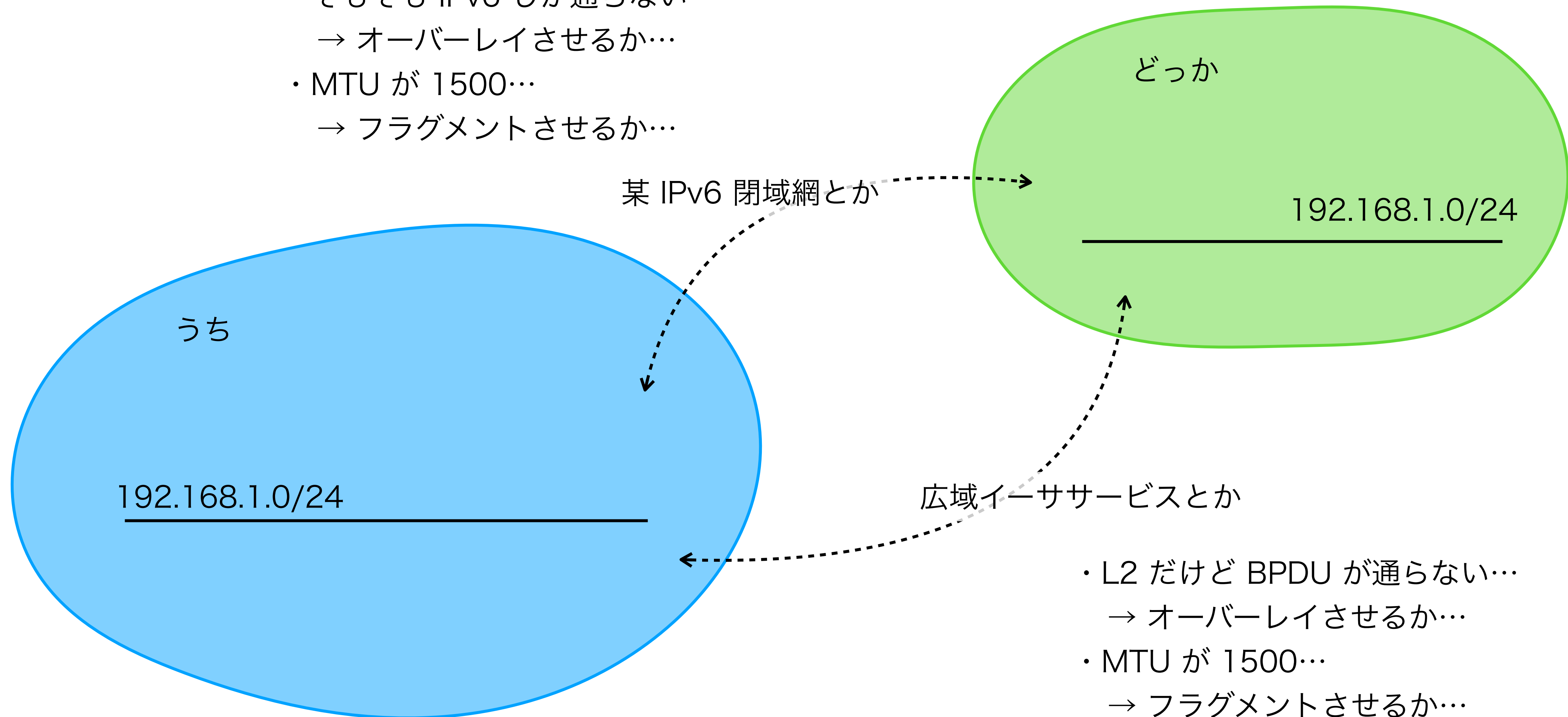
浅間正和@有限会社銀座堂

L2 を (MST で) のばしたい...



L2 を (MST で) のばしたい...

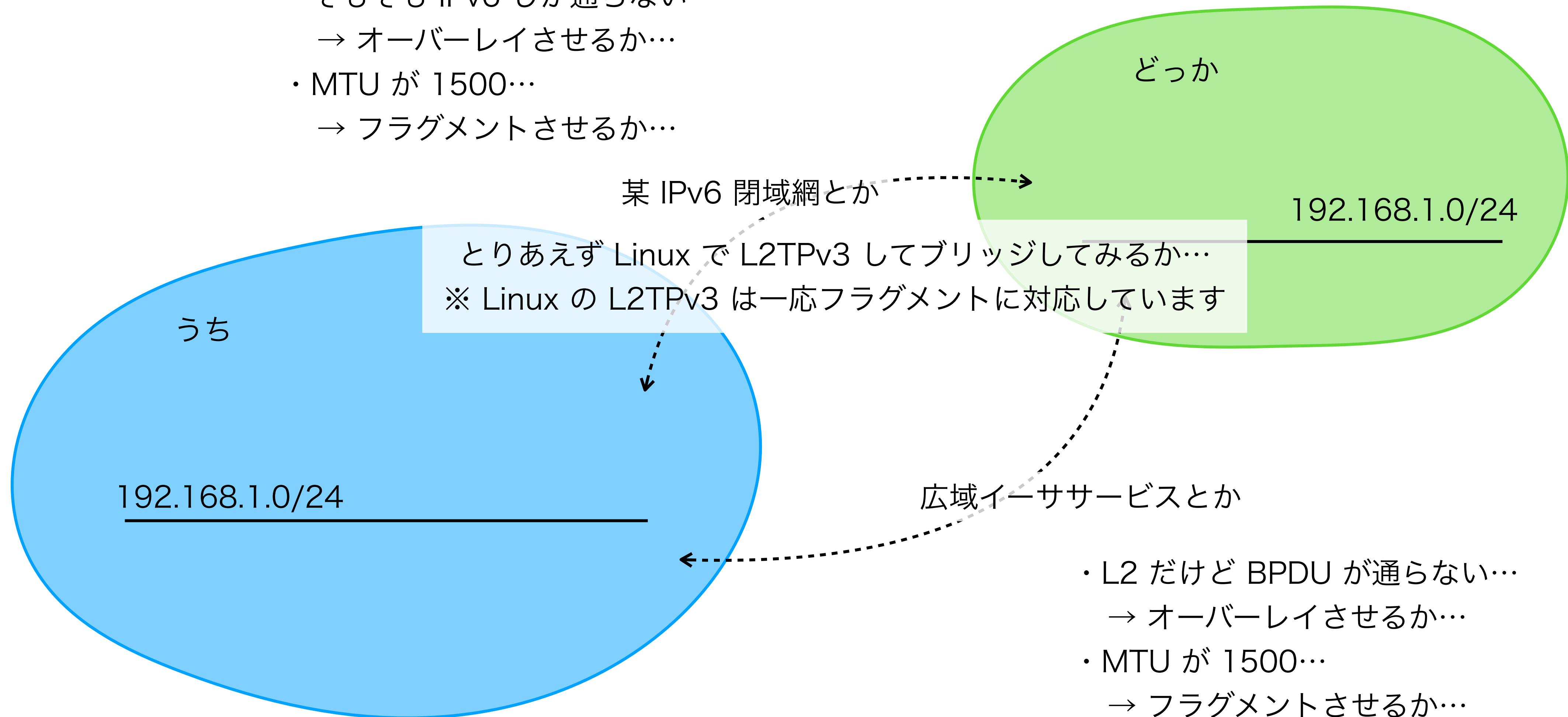
- そもそも IPv6 しか通らない...
 - オーバーレイさせるか...
- MTU が 1500...
 - フラグメントさせるか...



- L2 だけど BPDU が通らない...
 - オーバーレイさせるか...
- MTU が 1500...
 - フラグメントさせるか...

L2 を (MST で) のばしたい...

- そもそも IPv6 しか通らない...
 - オーバーレイさせるか...
- MTU が 1500...
 - フラグメントさせるか...

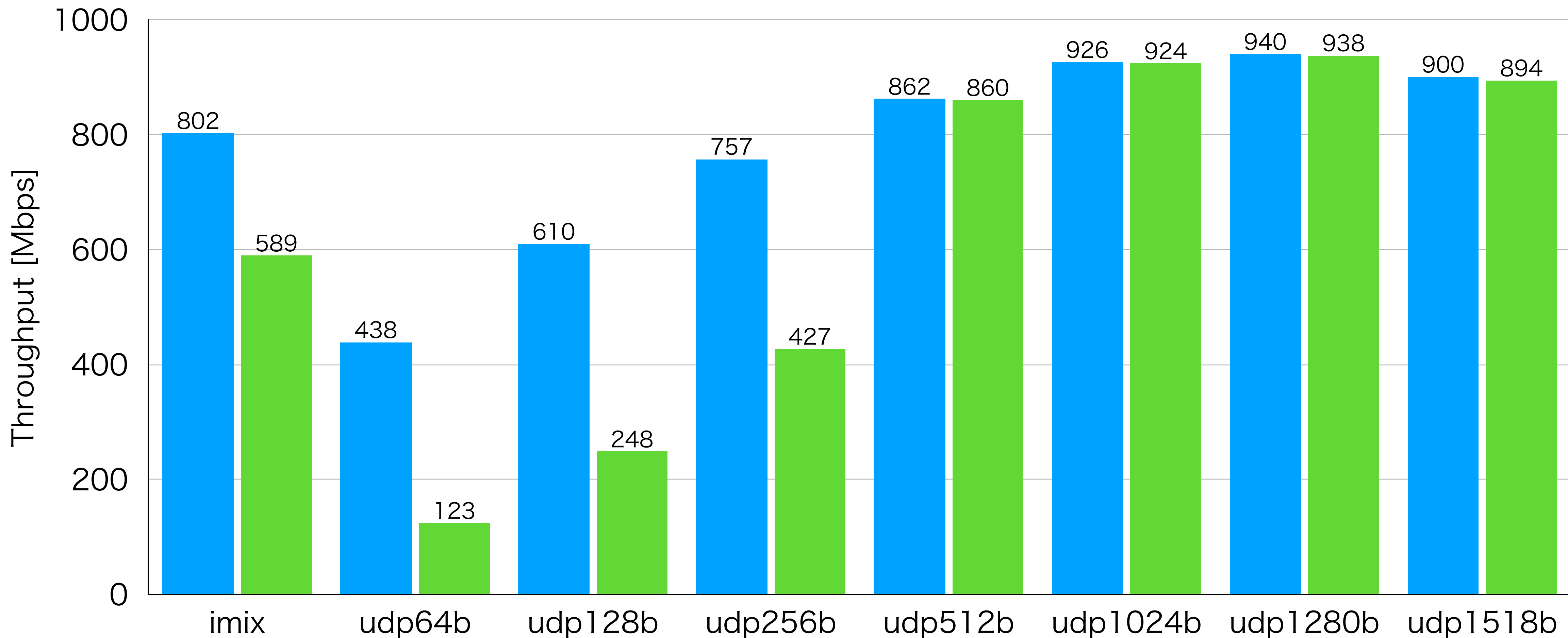


とりあえず Linux で L2TPv3 してブリッジしてみるか...
※ Linux の L2TPv3 は一応フラグメントに対応しています

- L2 だけど BPDU が通らない...
 - オーバーレイさせるか...
- MTU が 1500...
 - フラグメントさせるか...

■ wire-speed@l2tpv3/ipv6 ■ linux-l2tpv3/ipv6-bridge

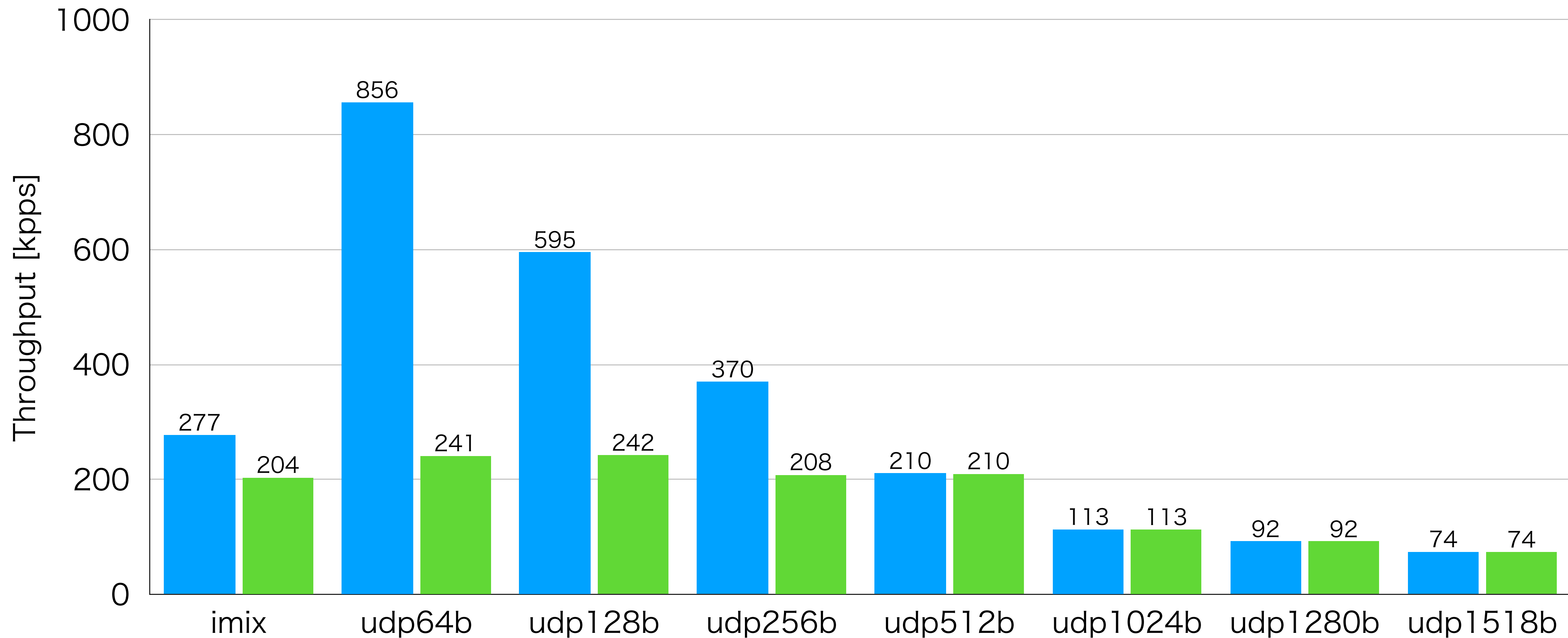
Skynew IN-1 ← Partaker C4 ← これらは何者かは後述



wire-speed@l2tpv3/ipv6

linux-l2tpv3/ipv6-bridge

Skynew IN-1 ← Partaker C4 ← これらは何者かは後述



ということで DPDK で書いてみた

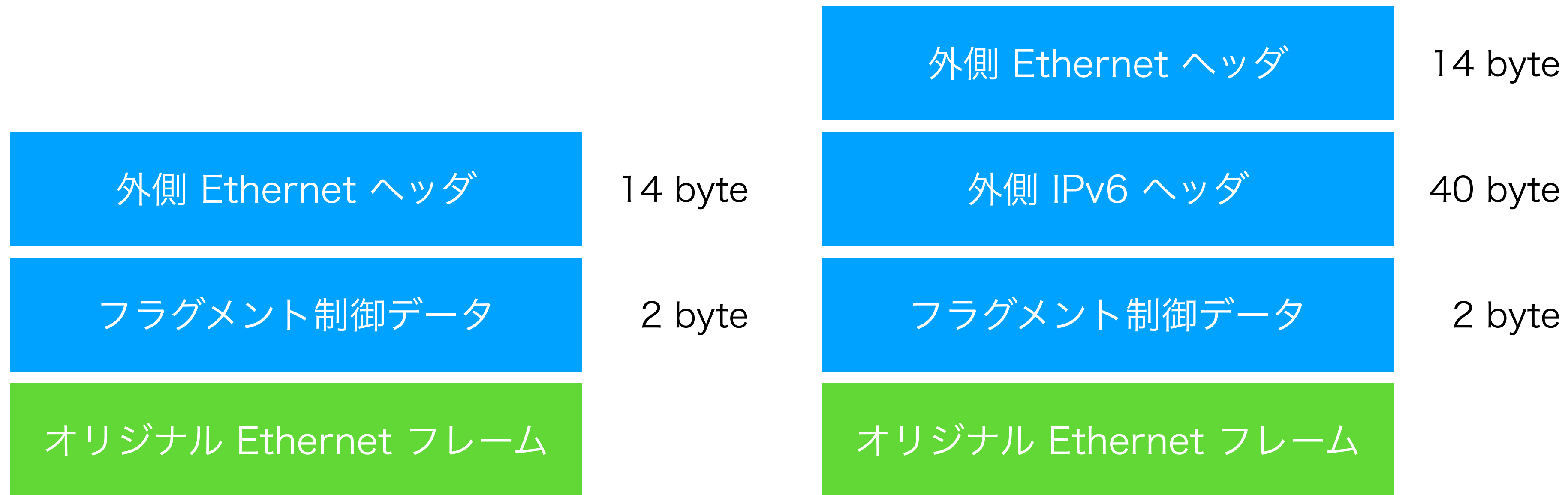
- Ethernet over Ethernet と Ethernet over IPv6 の両方に対応させる
 - ちょっとでもオーバーヘッドを抑えたい…
- over Ethernet でもフラグメントに対応させるために独自ヘッダとして 2 バイト使う (over IPv6 でもフラグメントは同じ仕組みにする)
 - L2TPv3 の 8 バイトよりは小さいのでそれよりは効率的なはず…
- リアセンブルはちゃんとするのが面倒なのでリオーダーはせずインオーダーで届くことを前提してしまいしかも 2 分割までしか対応しない
 - 完全同一な Ethernet/IPv6 ヘッダのパケットで入れ替わることなんてないですよ…
- over IPv6 は RS 出すのもめんどいので RA をひたすら待つ

おおまかな処理内容

- エンキャップ処理専用スレッドとデキャップ処理専用スレッドとそれ以外の処理スレッド(メインスレッド)の合計 3 つのスレッドをぶん回す
- エンキャップ処理専用スレッド:
 - DL(ダウンリンク)ポートから受信したパケットを受け取りエンキャップしたものを UL(アップリンク)ポートから送信するをひたすら繰り返す
 - UL ポートの MTU に収まらない時はフラグメントし送信する(後述)
- デキャップ処理専用スレッド:
 - UL ポートから受信したパケットを受け取りデキャップしたものを DL ポートから送信するをひたすら繰り返す
 - フラグメントされたパケットの場合はリアセンブルし送信する(後述)

エンキャップとデキャップ

- エンキャップ時のヘッダはこんな感じ



フラグメント制御データ

- フラグメントされていない場合:

0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0

- フラグメントされた前半の場合:

1 0 0 0 0 0 0 0 0 0 0 1 1 0 1 0 1

- フラグメントされた後半の場合:

1 1 0 0 0 0 0 0 0 0 0 1 1 0 1 0 1

↑ ID(リアセンブル時の判断用; IP の ID と同じようなもの?)(フラグメントする毎にインクリメント)
↑ 後半フラグ(1: フラグメントされた後半; 0: フラグメントされた前半)(フラグメントフラグ=1の時のみ)
↑ フラグメントフラグ(1: フラグメントされている; 0: フラグメントされていない)

Skynew IN-1

Partaker C4

外観



価格

26,990円(税込)
(Skynew)

38,421円(税込)
(Amazon)

CPU

Intel Celeron N2940 @ 1.83GHz
最大 2.25GHz 4C/4T

Intel Celeron J3160 @ 1.60GHz
最大 2.24GHz 4C/4T

メモリ

DDR3 1333MHz 4GiB

DDR3 1600MHz 4GiB

NIC

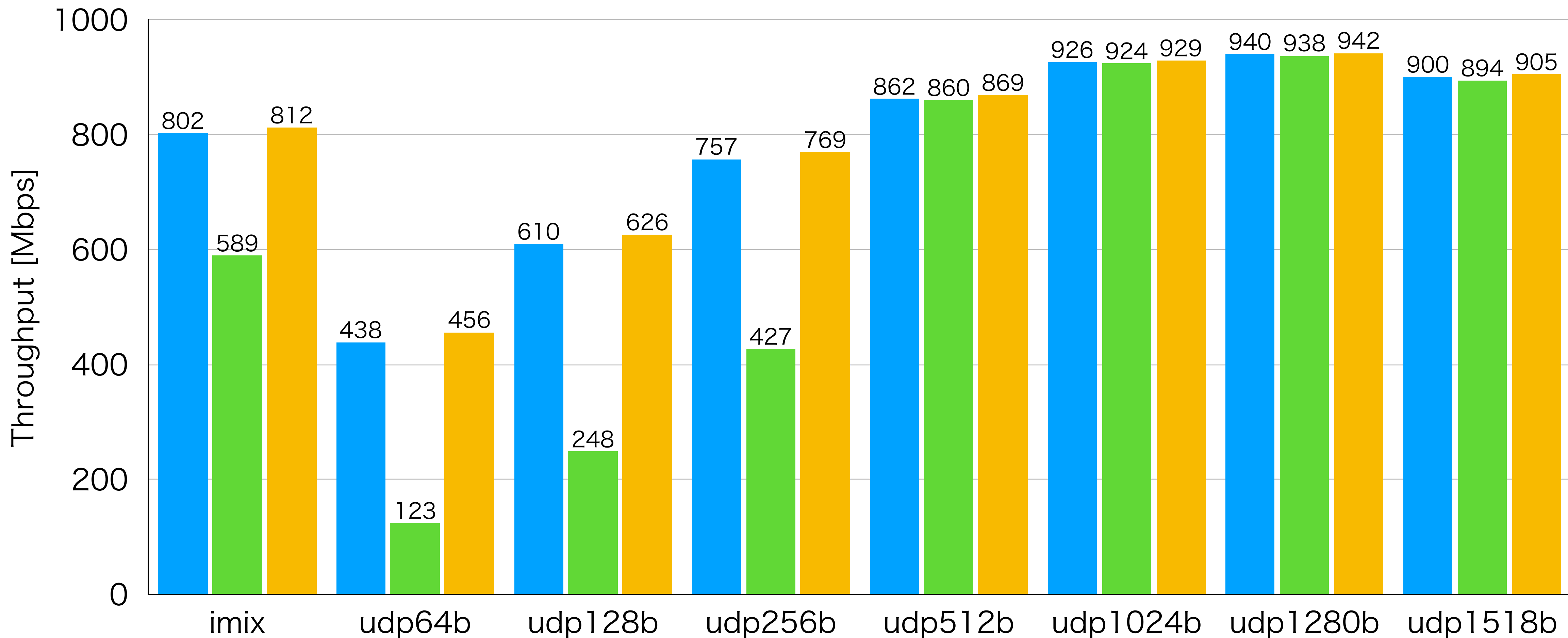
Intel I211 1000BASE-T x 4 ポート

Intel I211 1000BASE-T x 4 ポート

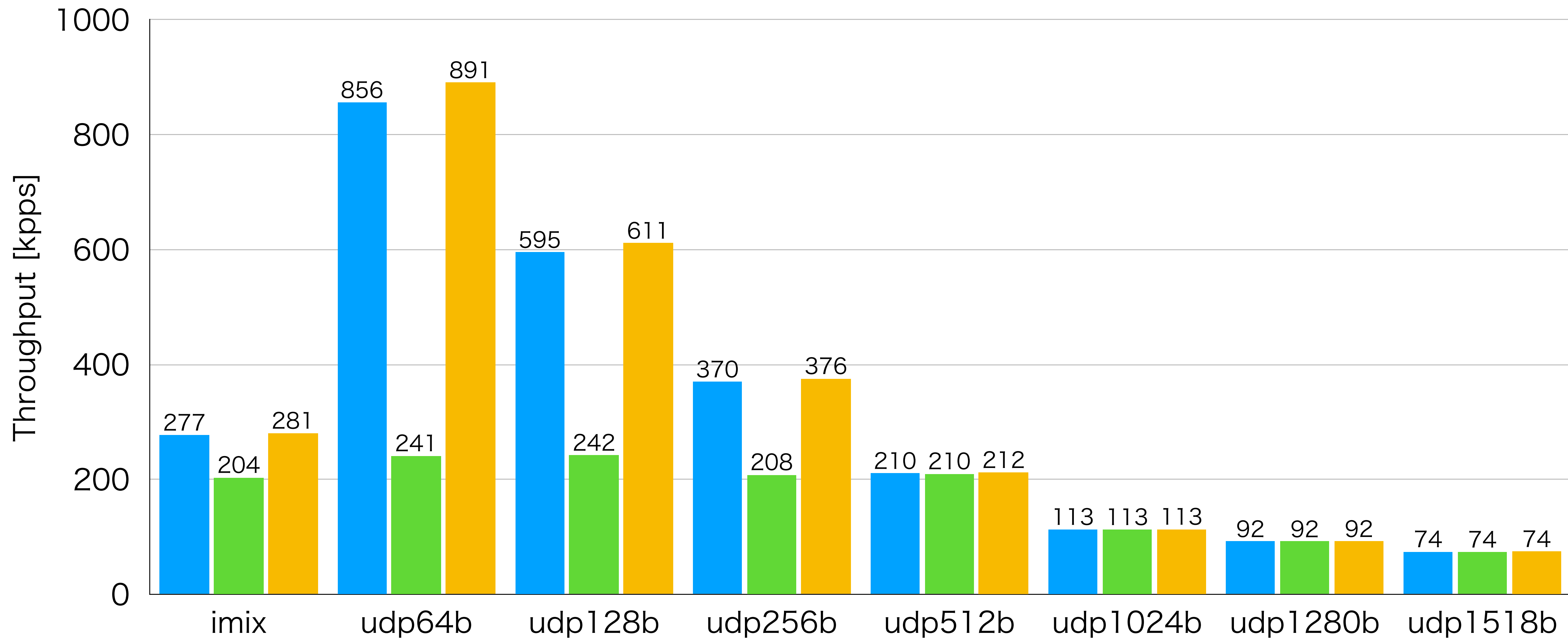
備考

Protectli Vault の OEM 品?

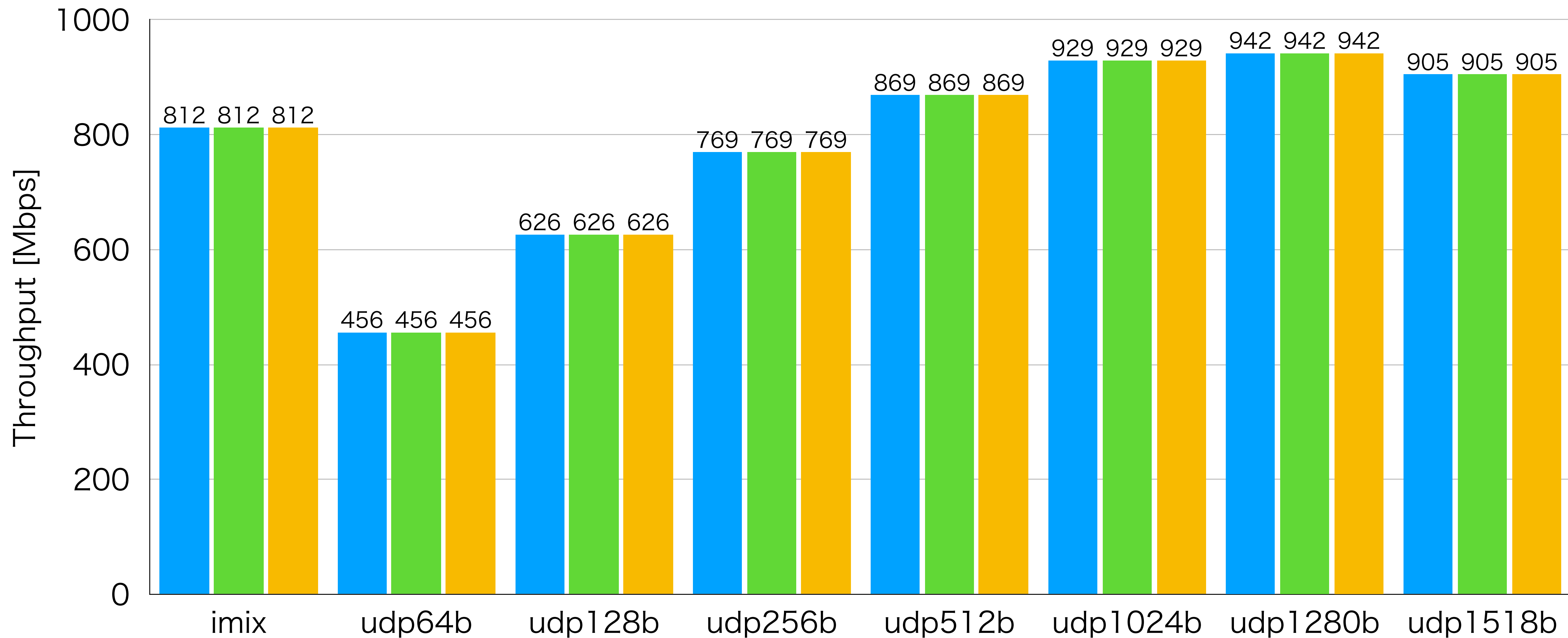
■ wire-speed@l2tpv3 ■ linux-l2tpv3-bridge ■ ginzado-pseudowire-ip6
Skynew IN-1 ← Partaker C4



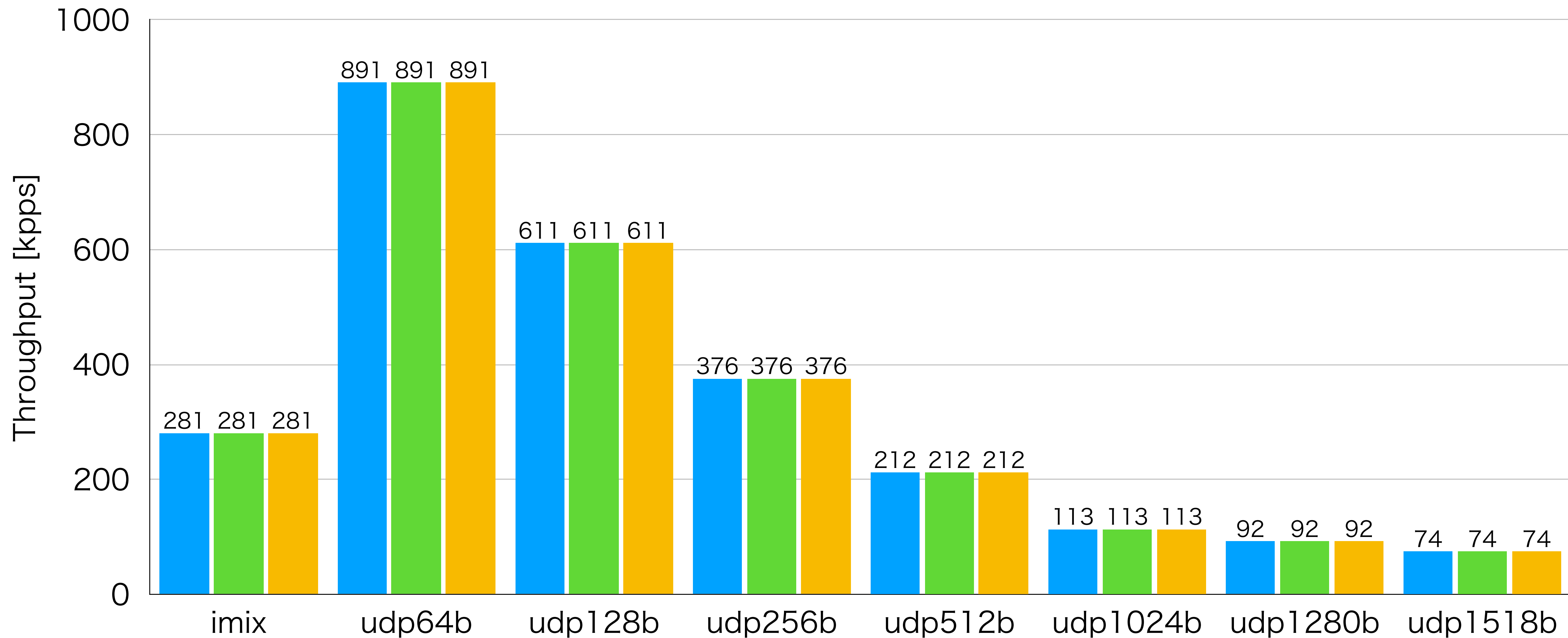
wire-speed@l2tpv3 linux-l2tpv3-bridge ginzado-pseudowire-ip6
Skynew IN-1 ← Partaker C4



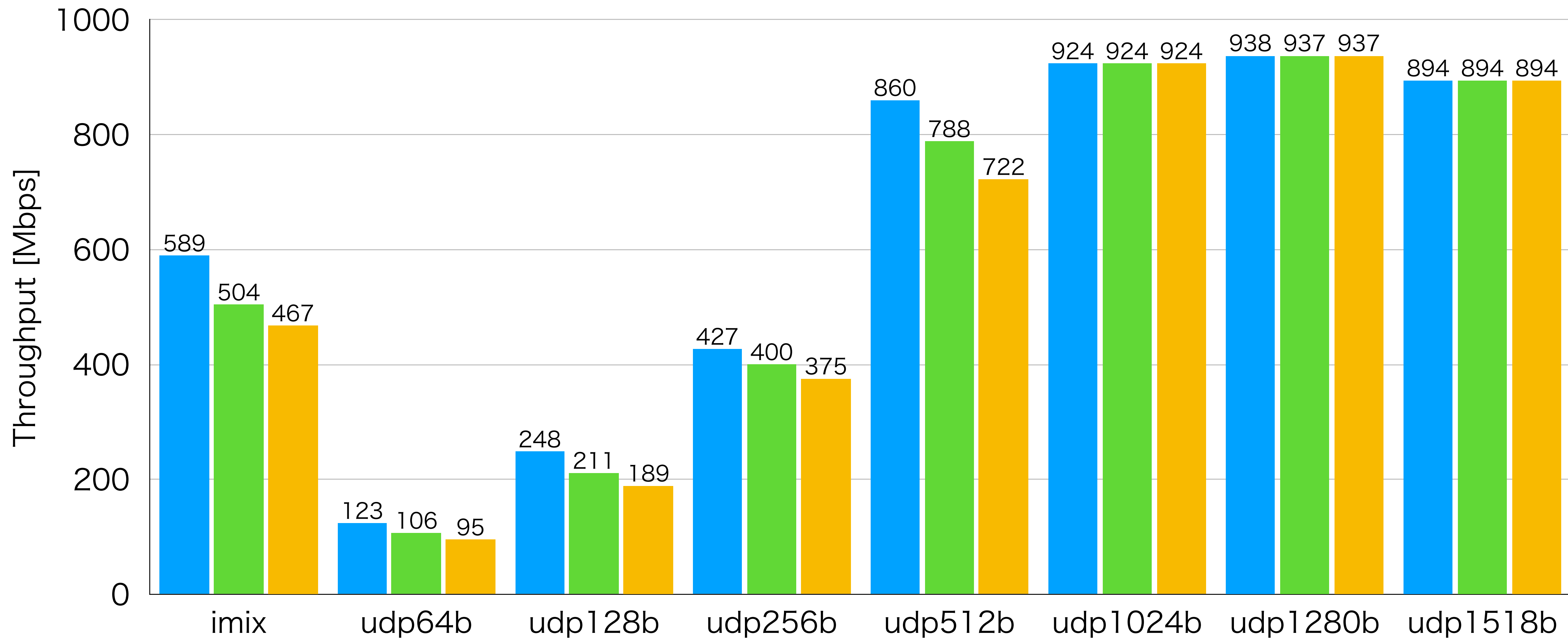
■ Skynew IN-1 ← Partaker C4 ■ Skynew IN-1 → Partaker C4 ■ Skynew IN-1 ⇔ Partaker C4
ginzado-pseudowire-ip6



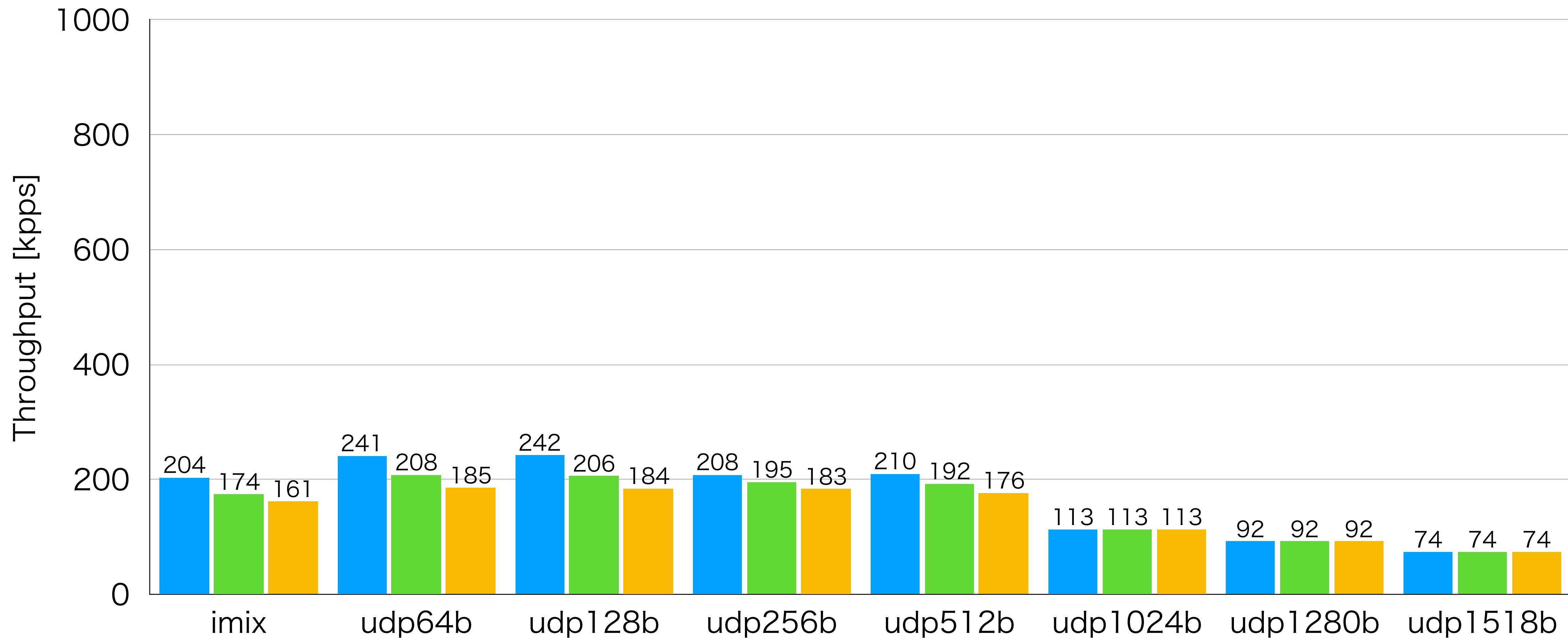
■ Skynew IN-1 ← Partaker C4 ■ Skynew IN-1 → Partaker C4 ■ Skynew IN-1 ⇌ Partaker C4
ginzado-pseudowire-ip6



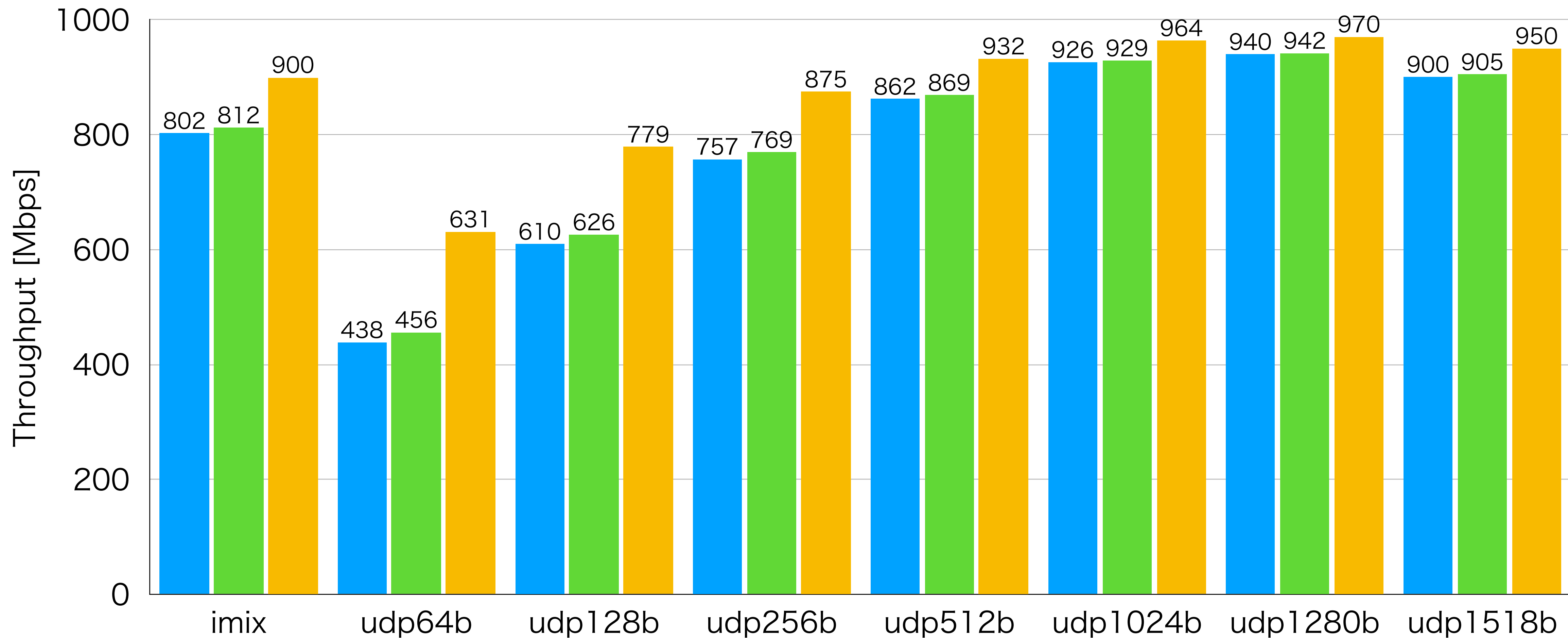
■ Skynew IN-1 ← Partaker C4 ■ Skynew IN-1 → Partaker C4 ■ Skynew IN-1 ⇌ Partaker C4
linux-l2tpv3/ipv6-bridge



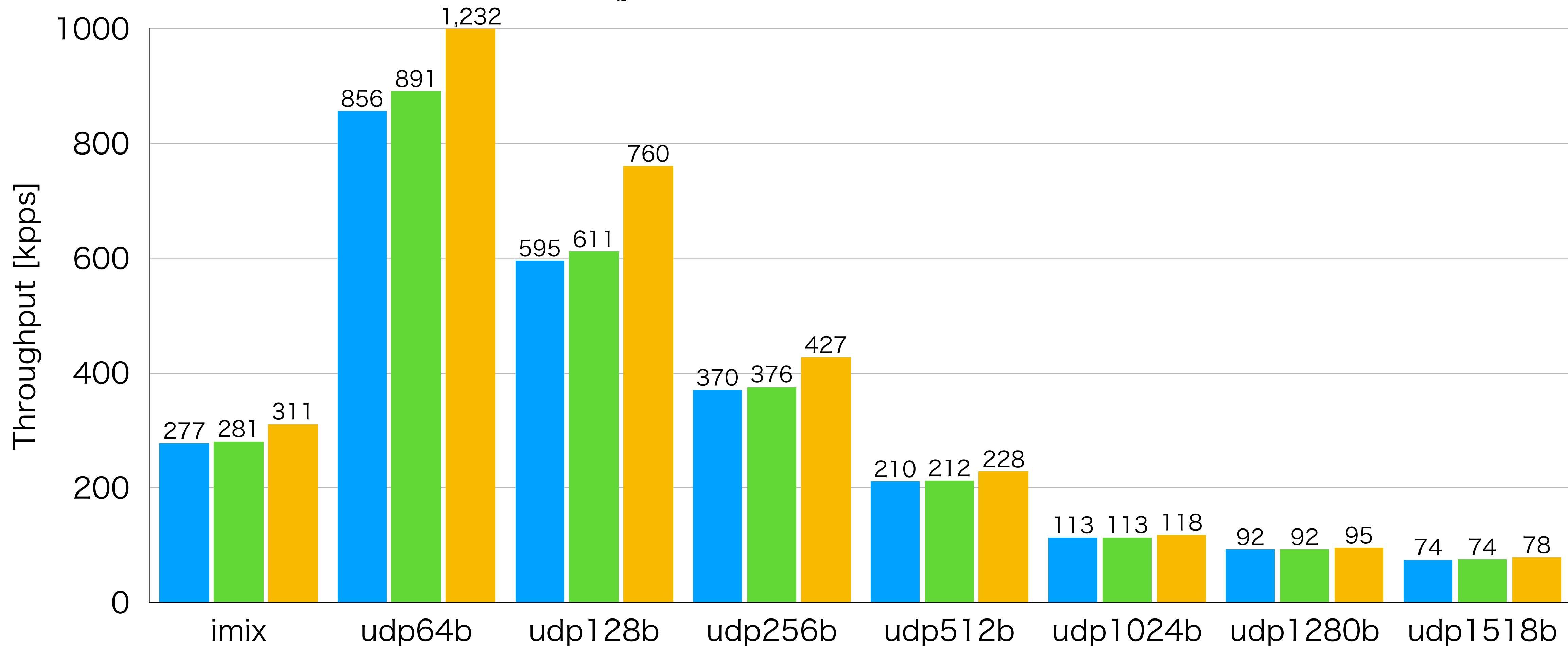
■ Skynew IN-1 ← Partaker C4 ■ Skynew IN-1 → Partaker C4 ■ Skynew IN-1 ⇔ Partaker C4
linux-l2tpv3/ipv6-bridge



■ wire-speed@l2tpv3 ■ ginzado-pseudowire-ip6 ■ ginzado-pseudowire-eth
Skynew IN-1 ← Partaker C4



wire-speed@l2tpv3 ginzado-pseudowire-ip6 ginzado-pseudowire-eth
Skynew IN-1 ← Partaker C4



まとめ

- 初期費用両端 60,000円(税込)くらい(某 IPv6 閉域網だけなら)維持費用も月額 12,000円(税込)くらいで L2 伸ばせます
- ただし電気代はもしかしたら結構くうかも？
 - 我が家ではこれを動かしている時に(因果関係は不明ですが) 2 回ブレーカーが落ちました
- ソースコード公開しました
 - <https://github.com/ginzado/dpdk>

おまけ: BPDU カウンタ

- MST でバックアップが生きてるのか死んでるのか監視したい
 - メインが落ちた際に「実はバックアップも(メインよりも前に)死んでたんだけど気づかず切り替え失敗」みたいなことが起こるのは避けたい
 - 統計カウンタに UL ポートで受け取った(エンキャップされた) BPDU フレームをカウントする機能を実装

```
$ ./gpwstats
gpwstats.ul_rx_packets          2
gpwstats.ul_rx_bytes           120
gpwstats.ul_rx_bpdu            0
gpwstats.ul_tx_packets          2
gpwstats.ul_tx_bytes           152
... 省略 ...
```

- これを使って BPDU フレームが流れてこなくなったらアラートを上げる Nagios プラグインも作成(GitHub のリポジトリにも入ってます)