



さくらのクラウドにおける L2ネットワークの仕様について

～DSR構成のLVS構築時にお気をつけいただくために～



<https://www.sakura.ad.jp/>

DAY

2018/9

COMPANY

さくらインターネット株式会社

DEPARTMENT

クラウドチーム

NAME



サーバAからサーバBに通信をしようとしています。
この時点ではスイッチのMACアドレステーブルは空の状態とします。

MACアドレステーブル

ポート	MACアドレス



サーバBに
通信したい



IP : 192.168.0.1
MAC : AAAA



IP : 192.168.0.2
MAC : BBBB



IP : 192.168.0.3
MAC : CCCC



IP : 192.168.0.4
MAC : DDDD



TCP/IPのネットワークでは、通信を行う際、送信先の**IPアドレス**と**MACアドレス**が必要となります。そのアドレス解決を行うプロトコルを**ARP**と言います。

サーバAはサーバBのMACアドレスを知る必要があるため、そのリクエスト (ARPリクエスト) をブロードキャスト (宛先 FF:FF:FF:FF:FF:FF) で送信します。

MACアドレステーブル

ポート	MACアドレス



ARPリクエスト



192.168.0.2の機器のMACアドレスを知りたい



IP : 192.168.0.1
MAC : AAAA



IP : 192.168.0.2
MAC : BBBB



IP : 192.168.0.3
MAC : CCCC



IP : 192.168.0.4
MAC : DDDD

arpテーブル

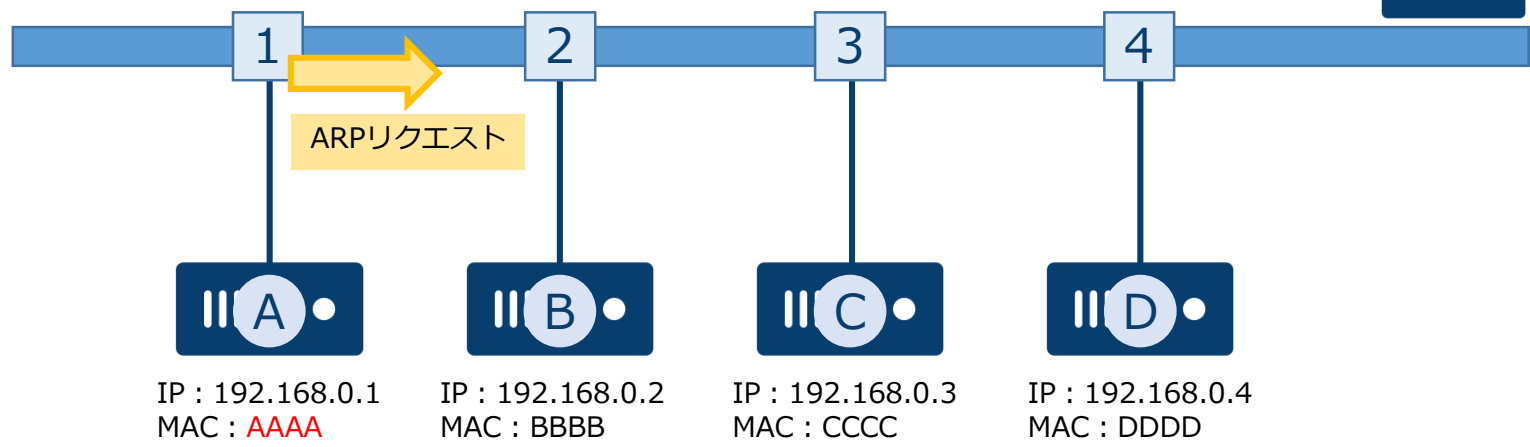
IPアドレス	MACアドレス



スイッチのポート1でARPリクエストを受信した時点で、MACアドレステーブルに送信元（サーバA）の情報が登録されます。

MACアドレステーブル

ポート	MACアドレス
ポート1	AAAA



arpテーブル

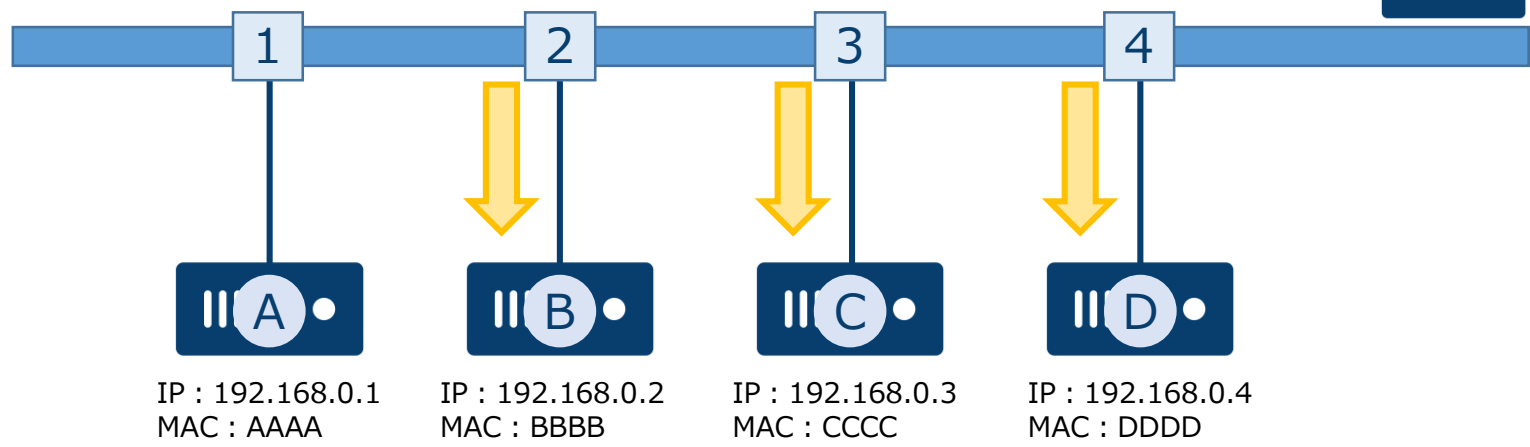
IPアドレス	MACアドレス
--------	---------



ARPリクエストが同一セグメントにブロードキャストされます。

MACアドレステーブル

ポート	MACアドレス
ポート1	AAAA



arpテーブル

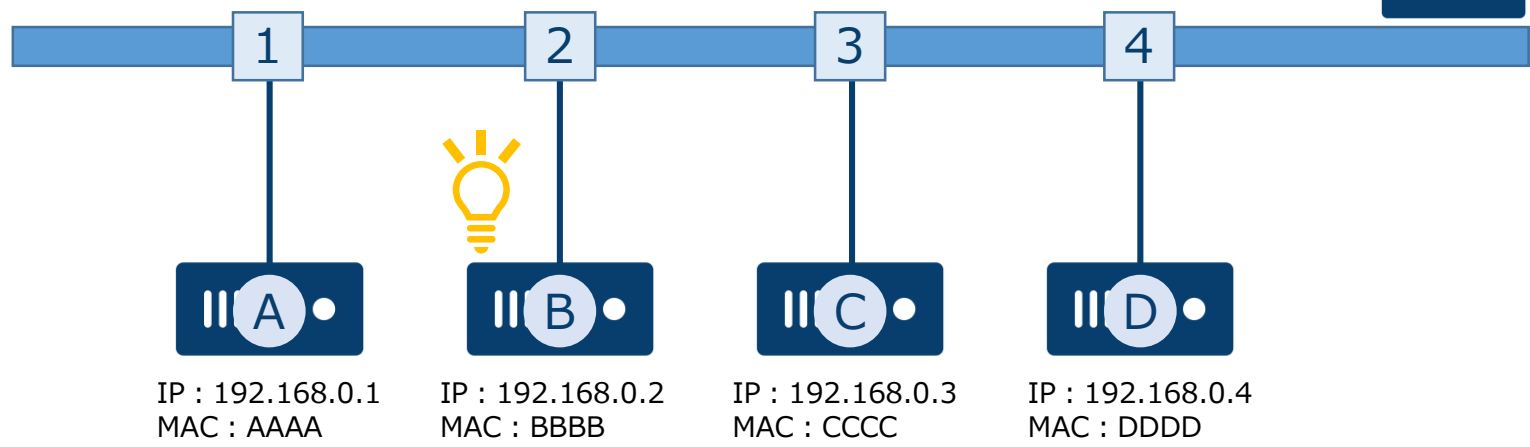
IPアドレス	MACアドレス



サーバBはARPリクエストの中身を見て、自分宛のリクエストだと気付きます。

MACアドレステーブル

ポート	MACアドレス
ポート1	AAAA



arpテーブル

IPアドレス	MACアドレス

192.168.0.2の機器は自分！

自分ではないので無視

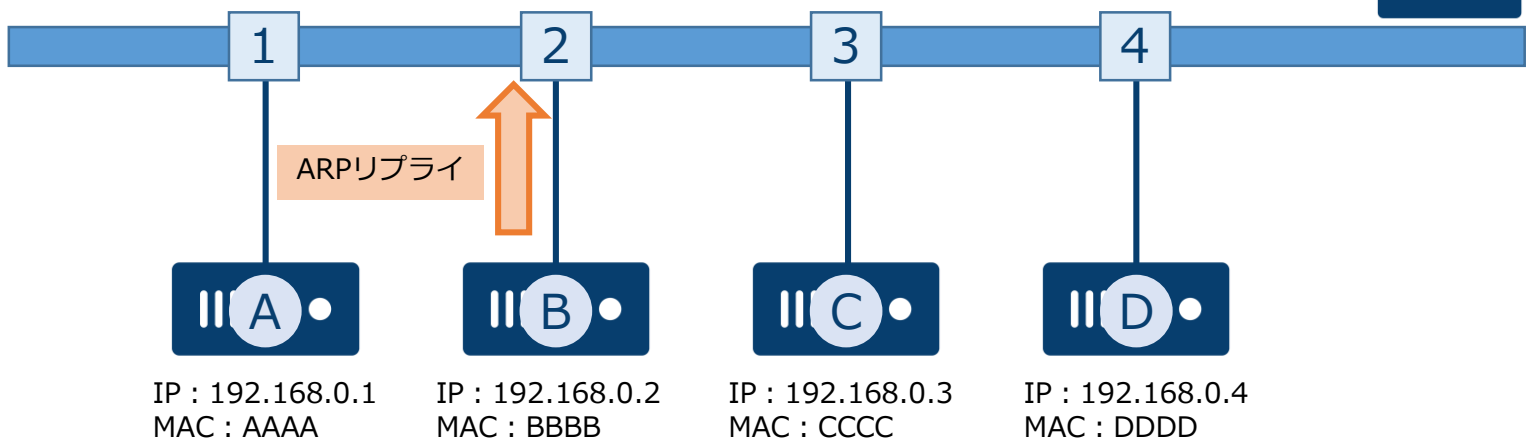
自分ではないので無視



サーバBは自身のMACアドレスをサーバAに伝えるため、ARPリプライを送信します。
 ※送信先はサーバAのみ

MACアドレステーブル

ポート	MACアドレス
ポート1	AAAA



arpテーブル

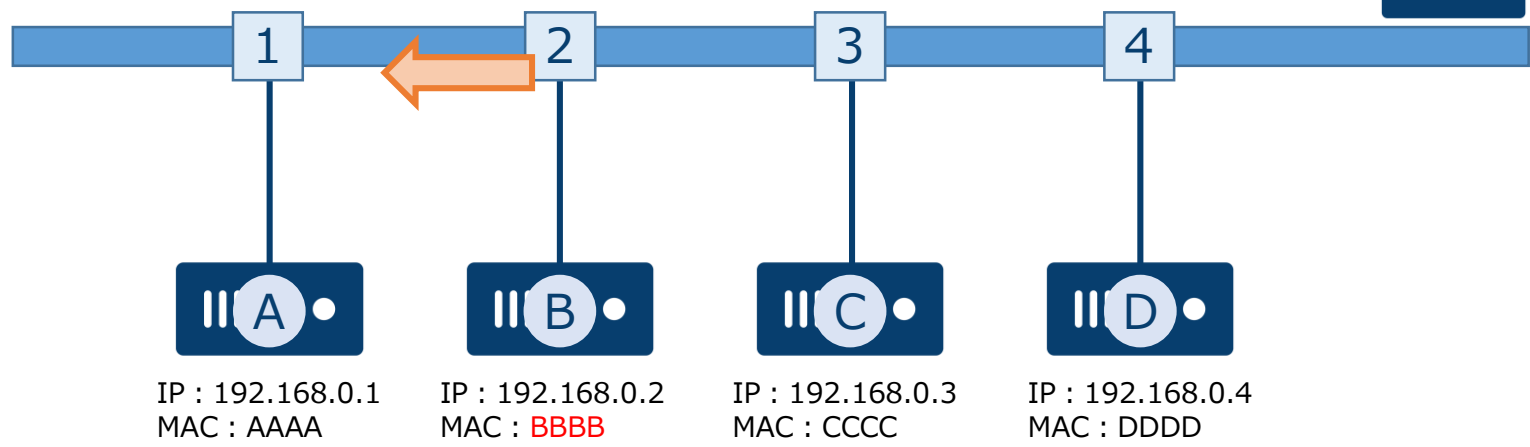
IPアドレス	MACアドレス



スイッチのポート2でARPリプライを受信した時点で、MACアドレステーブルに送信元（サーバB）の情報が登録されます。

MACアドレステーブル

ポート	MACアドレス
ポート1	AAAA
ポート2	BBBB



arpテーブル

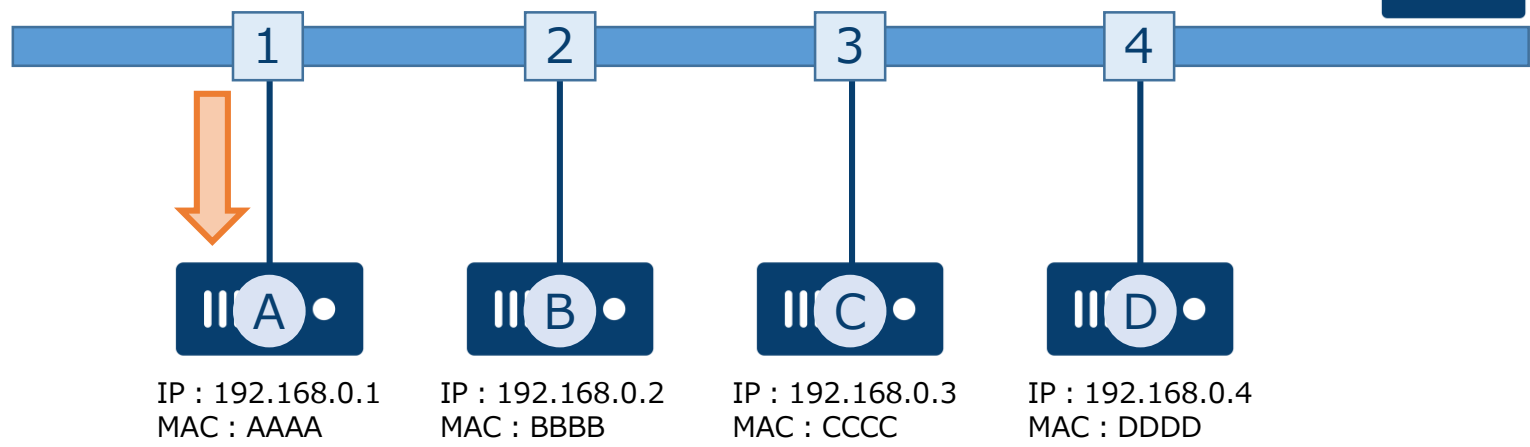
IPアドレス	MACアドレス



ARPリプライがサーバAに届き、サーバAはサーバBのMACアドレスを知ることができました。

MACアドレステーブル

ポート	MACアドレス
ポート1	AAAA
ポート2	BBBB

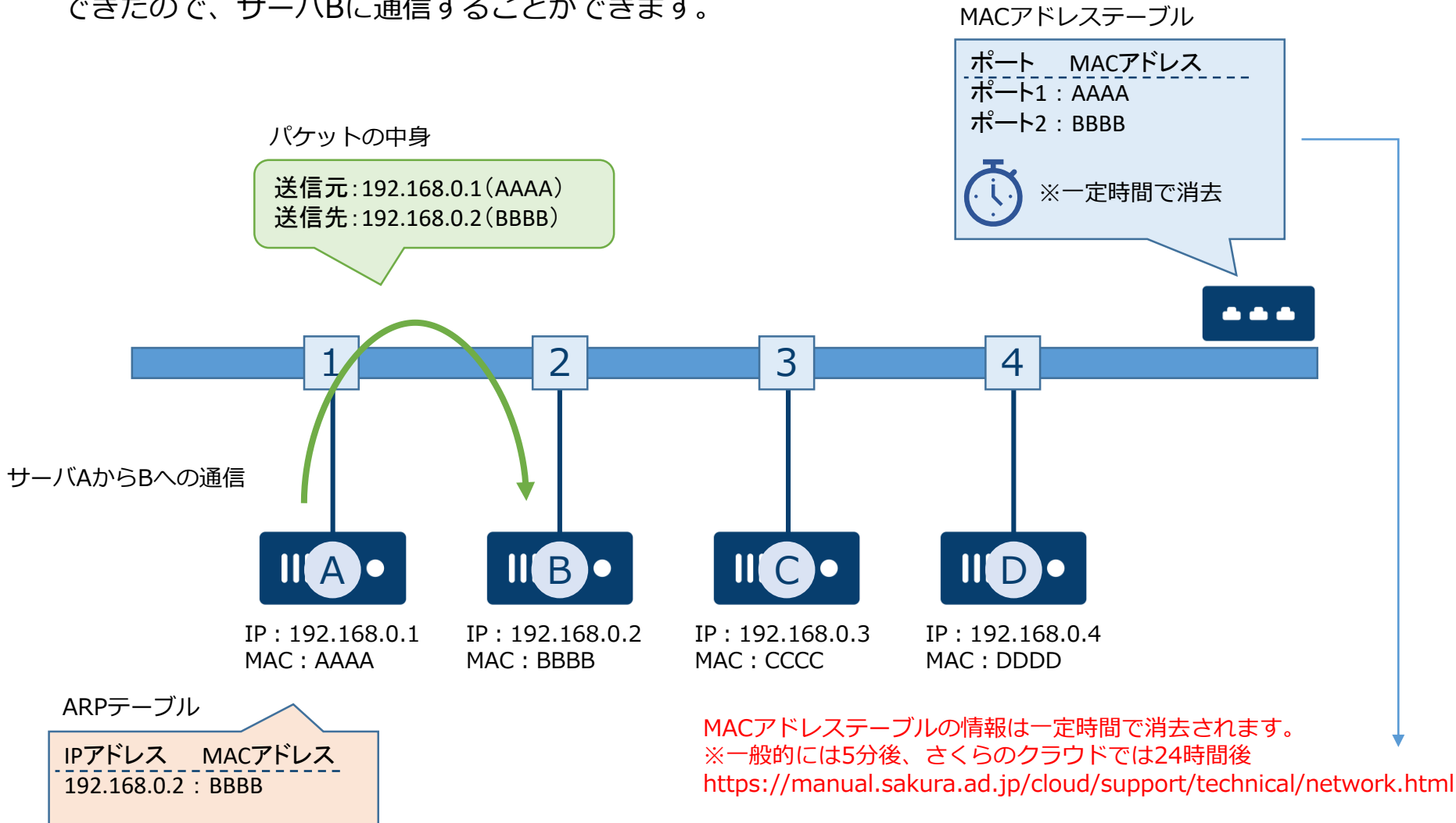


ARPテーブル

IPアドレス	MACアドレス
192.168.0.2	BBBB



サーバAはサーバB (送信先) のMACアドレスを知ることができたので、サーバBに通信することができます。

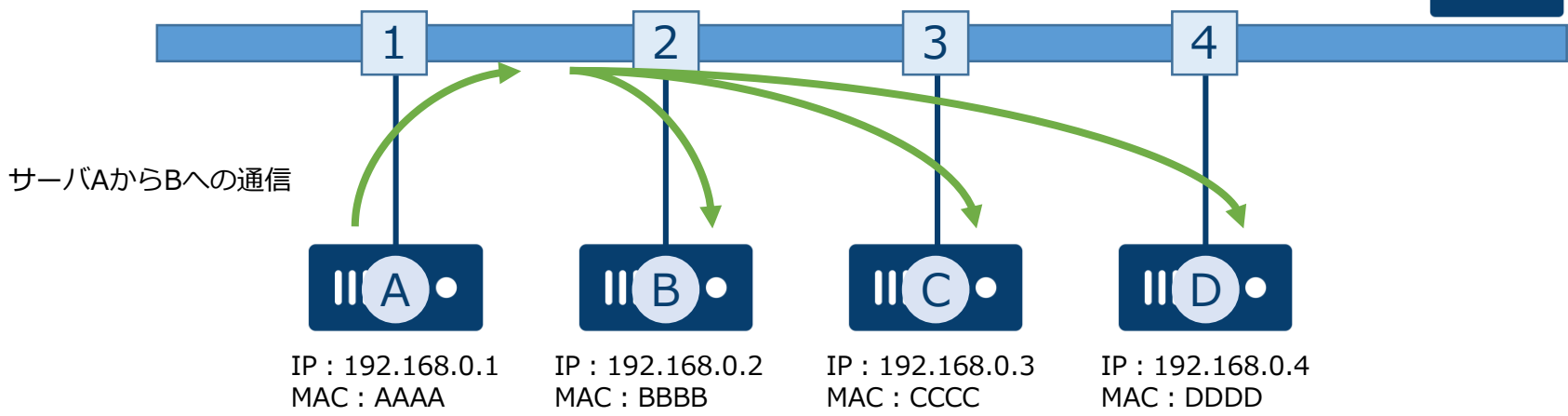


- サーバAからBへの通信時、
- ・スイッチのMACアドレステーブルは空
 - ・サーバAのARPテーブルに情報あり

この場合、サーバAからの通信は送信元以外のポートにフラッディングされます。

MACアドレステーブル

ポート	MACアドレス
空	



ARPテーブル

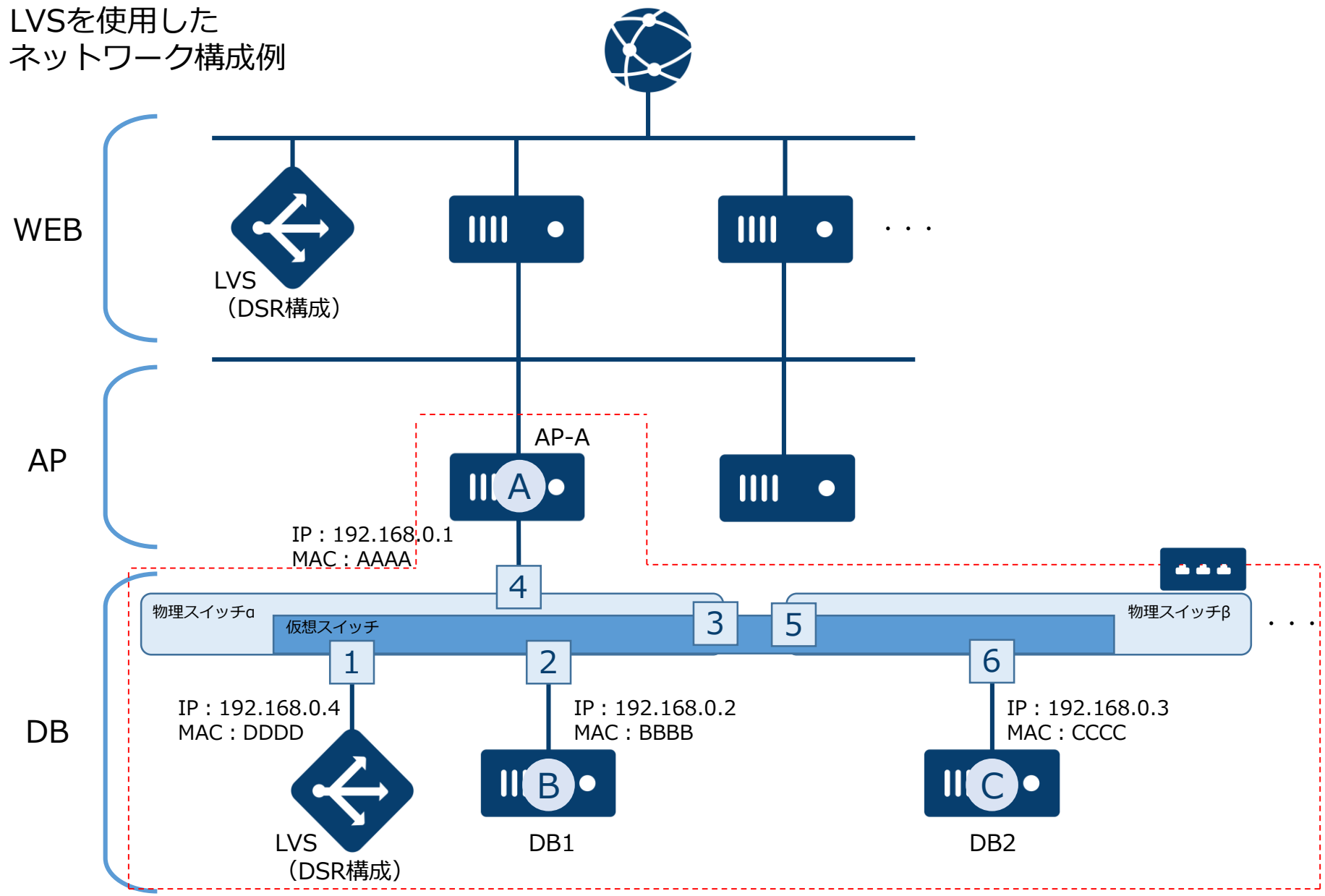
IPアドレス	MACアドレス
192.168.0.2	BBBB

フラッディングの通信（1対多）が発生した場合、回線が輻輳し他のお客様へ影響を与える可能性があります。そのため、さくらのクラウドではフラッディングパケットの総量帯域が制限されています。
<https://manual.sakura.ad.jp/cloud/support/technical/network-settings.html#dsr>



DSR構成のLVSを構築するケース

LVSを使用した
ネットワーク構成例

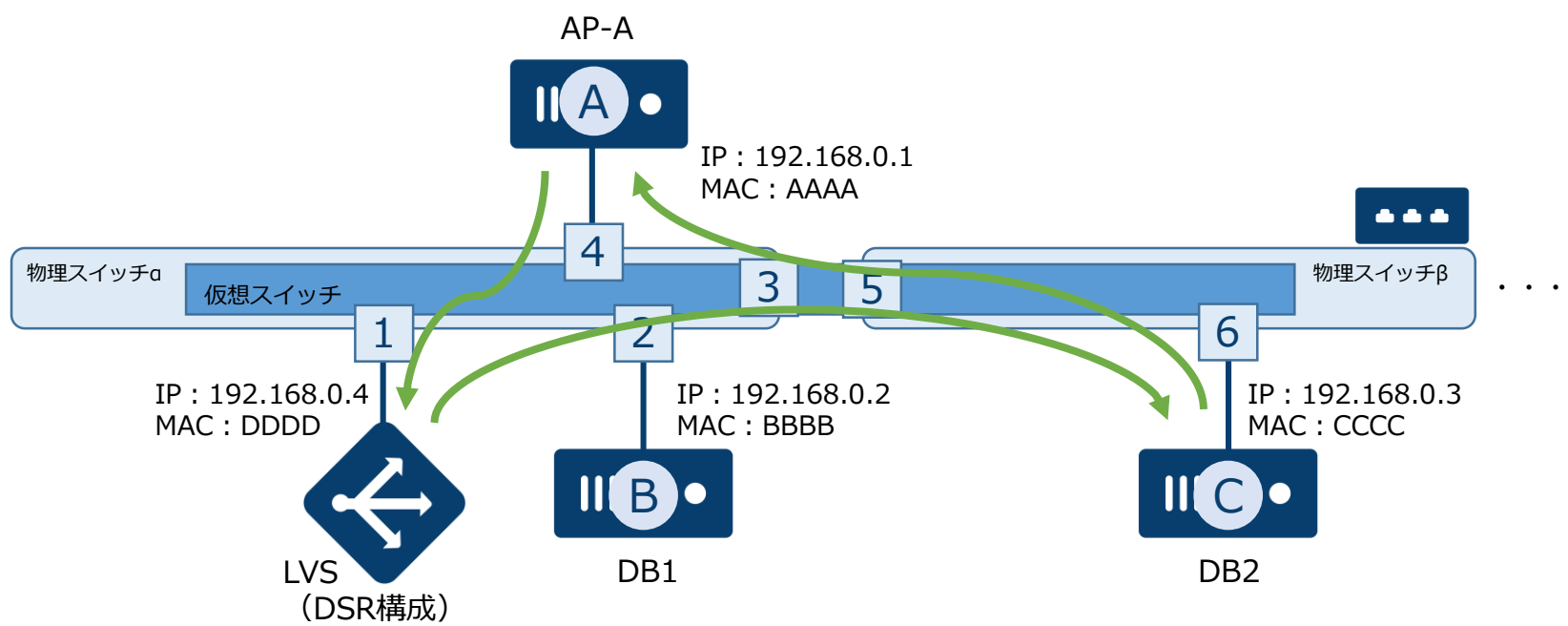




APサーバからDBへの通信をLVSで分散する部分を抽出した図になります。

なお、今回はDSR構成ですので、DBサーバに分散された通信は、LVSを経由せず直接APサーバに戻ります。

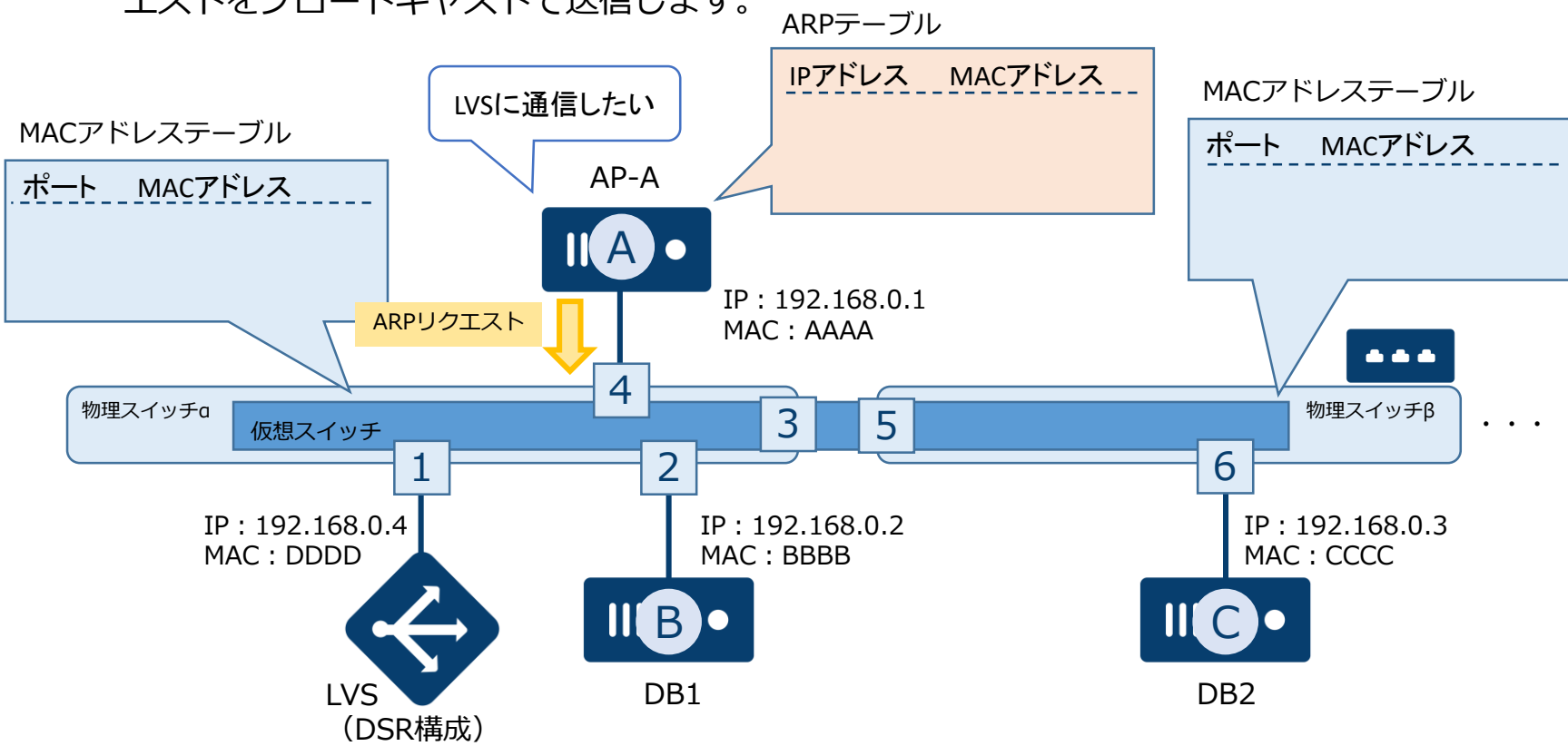
※以下は、LVSがサーバCを分散先としたケース





通信初期の段階では、サーバAのARPテーブル、各物理スイッチのMACアドレステーブルは空の状態とします。

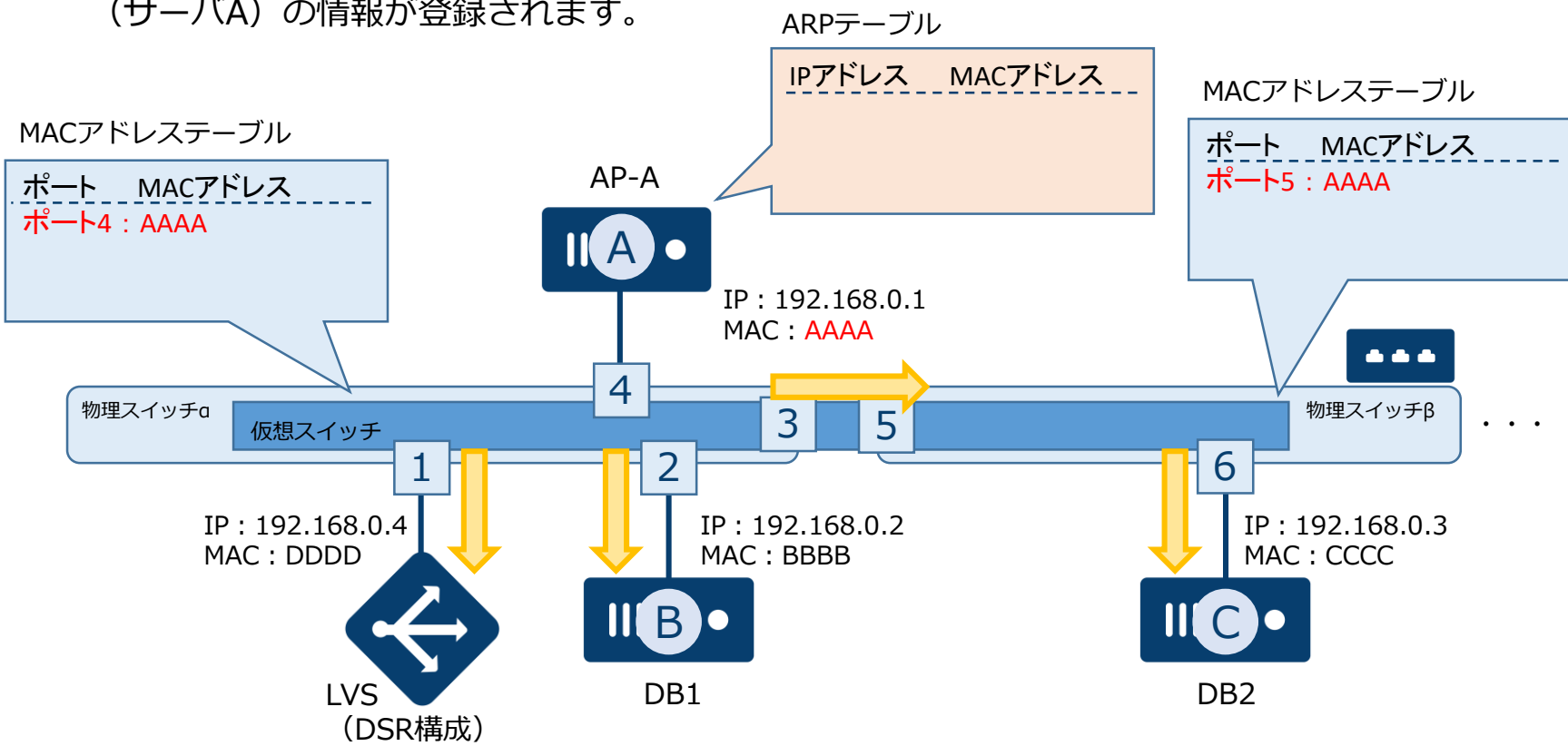
まず、サーバAはLVSに通信するためにARPリクエストをブロードキャストで送信します。





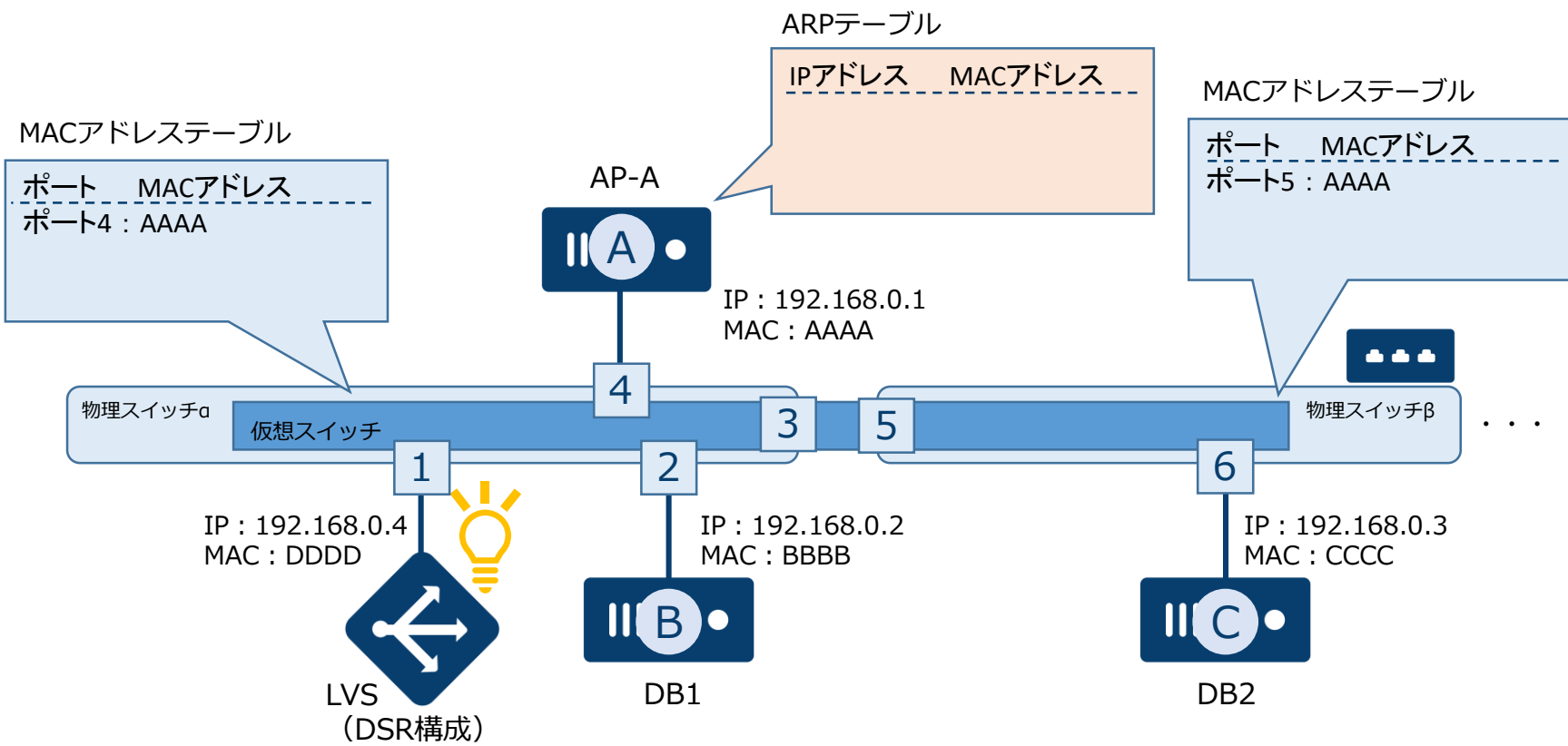
ARPリクエストが同一セグメントにブロードキャストされます。

スイッチの各ポートでARPリクエストを受信した時点で、MACアドレステーブルに送信元（サーバA）の情報が登録されます。



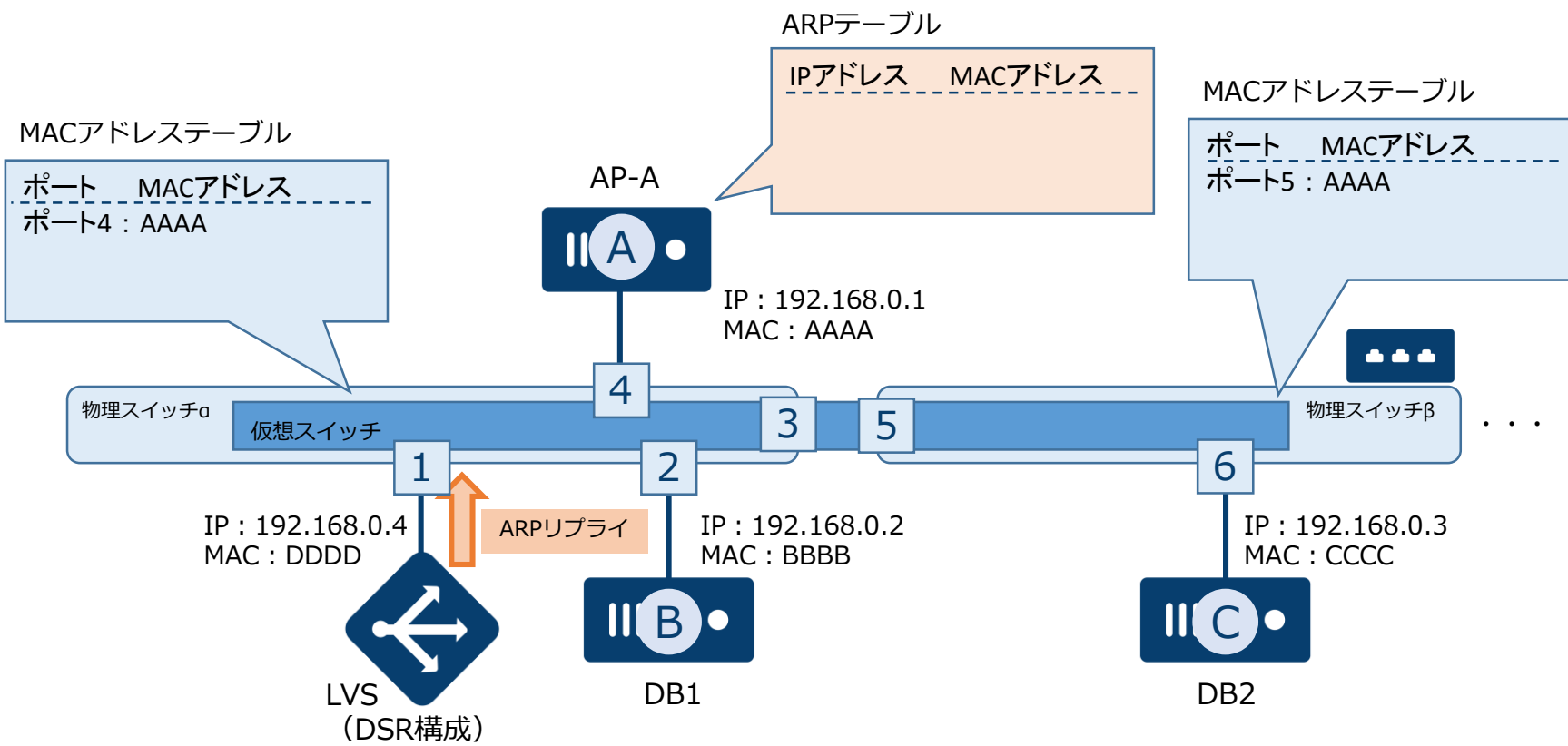


LVSはARPリクエストの中身を見て、自分宛のリクエストだと気づきます。



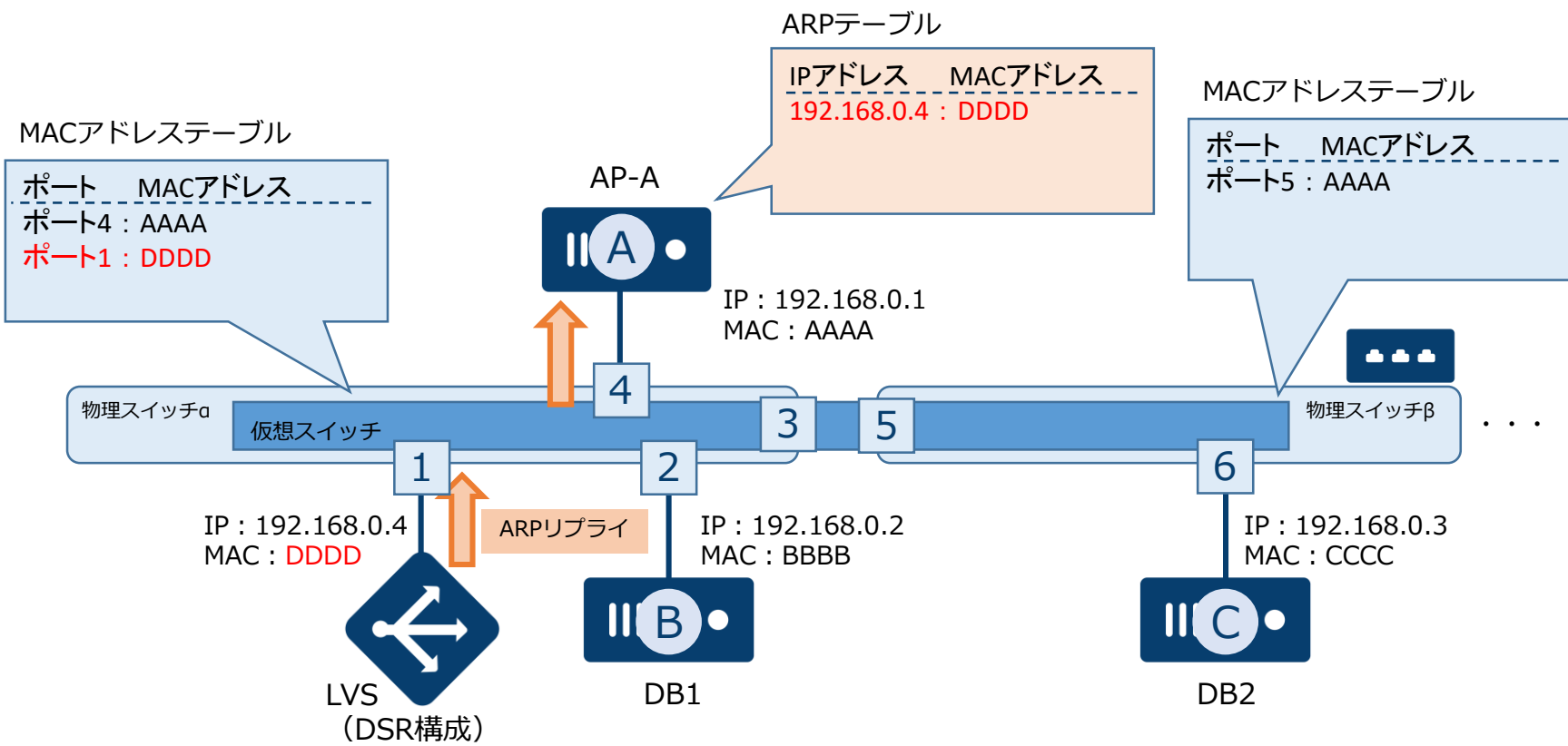


LVSは自身のサーバAにARPリプライを送信します。
 ※送信先はサーバAのみ





ARPリプライがサーバAに届き、サーバAはLVSのMACアドレスを知ることができました。





サーバAはLVISのMACアドレスを知ることができたので、LVISに通信することができます。

パケットの中身

送信元: 192.168.0.1 (AAAA)
 送信先: 192.168.0.4 (DDDD)

送信元 IP: クライアント
 送信元 Mac: クライアント
 あて先 IP: LVIS (VIP)
 あて先 Mac: 振り分け先サーバ

ARPテーブル

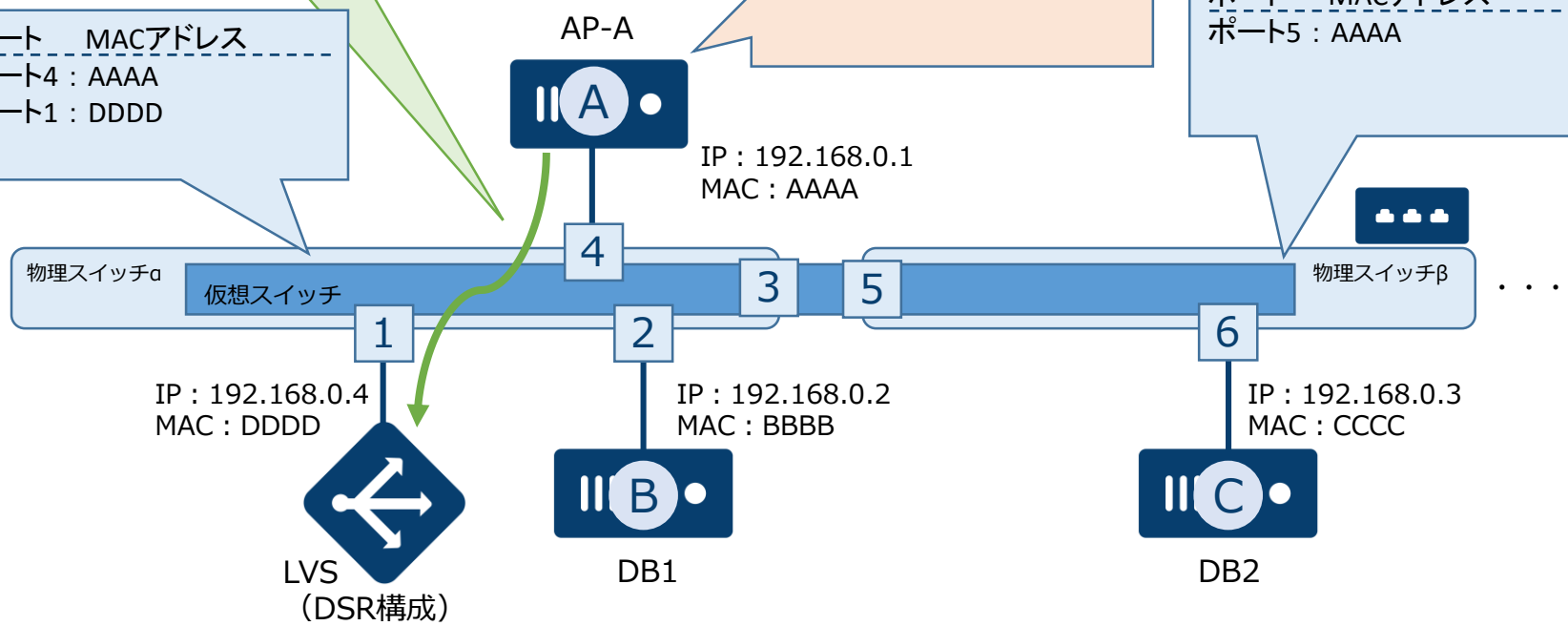
IPアドレス	MACアドレス
192.168.0.4	DDDD

MACアドレステーブル

ポート	MACアドレス
ポート4	AAAA
ポート1	DDDD

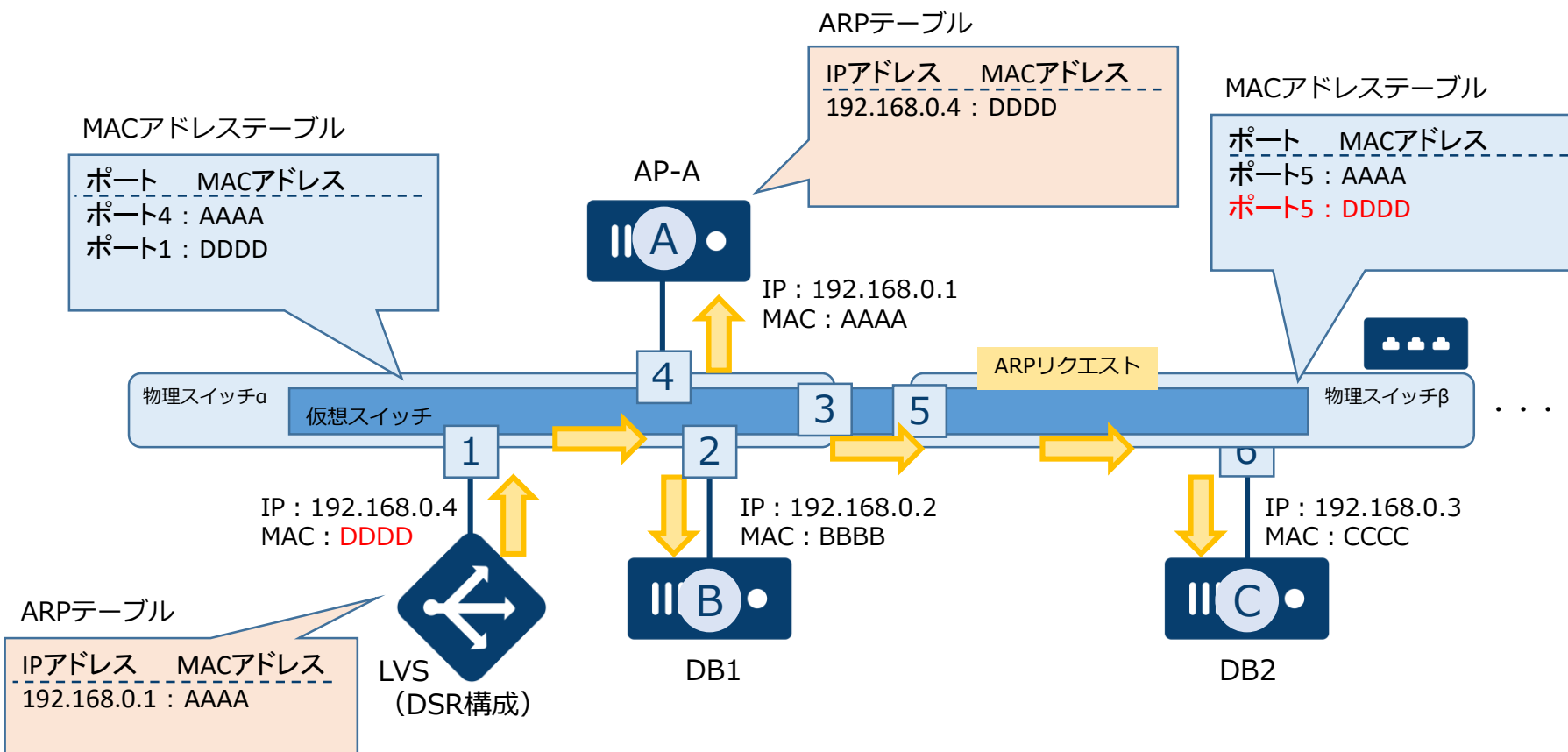
MACアドレステーブル

ポート	MACアドレス
ポート5	AAAA



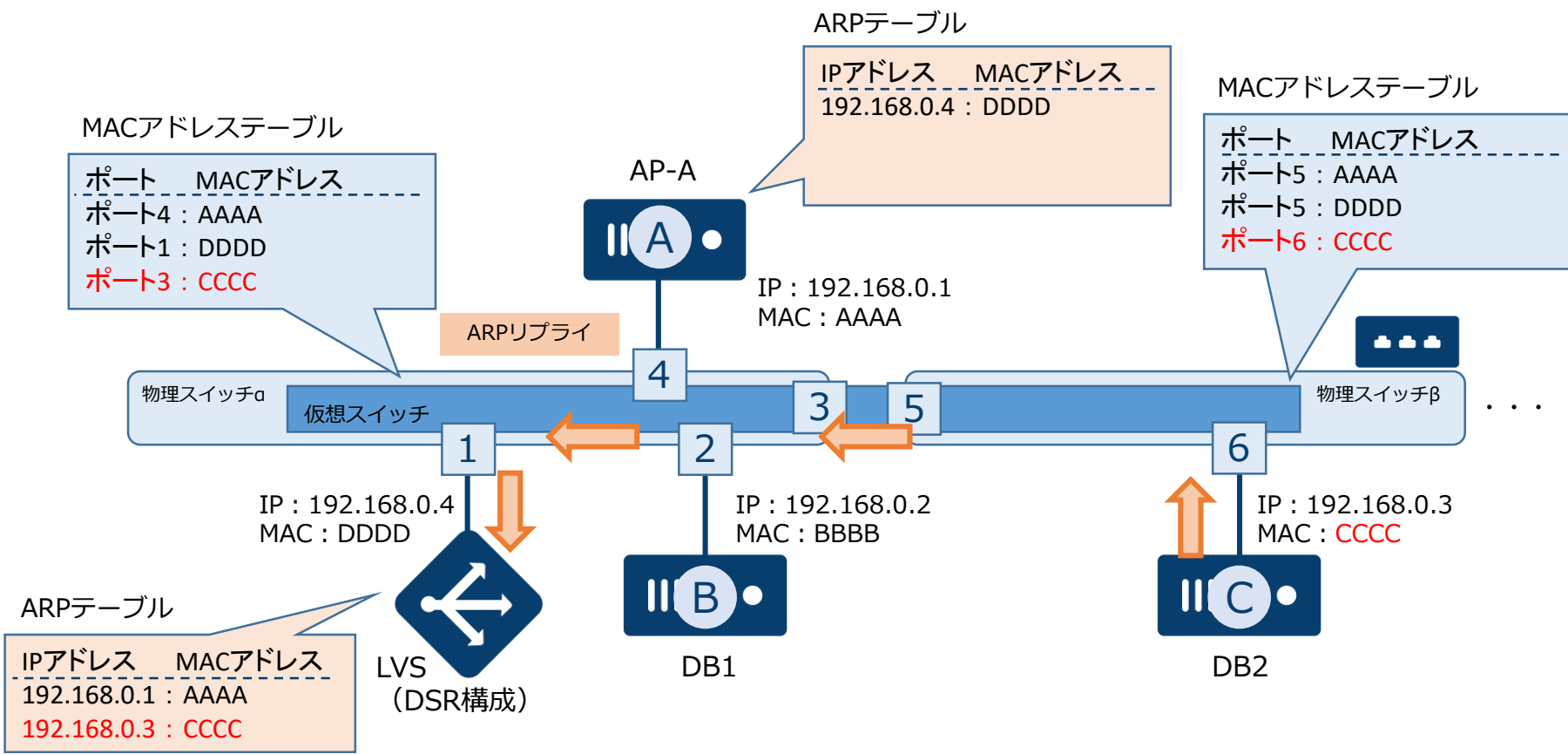


LVSは分散先のサーバとしてサーバCを選択したとします。
 LVSはサーバCのMACアドレスを調べるため、ARPリクエストを送信します。



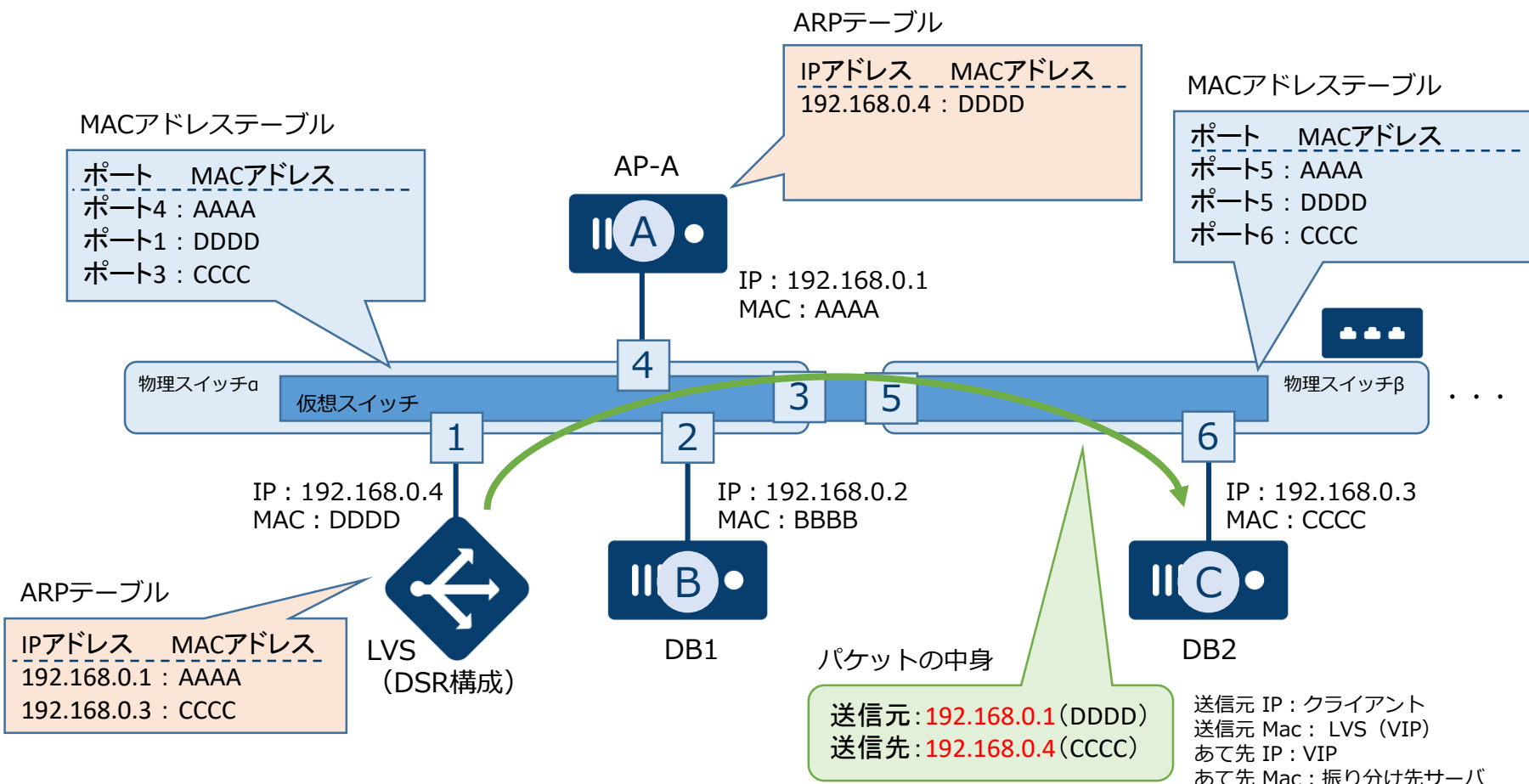


ARPリプライが届き、LVSはサーバCのMACアドレスを知ることができました。





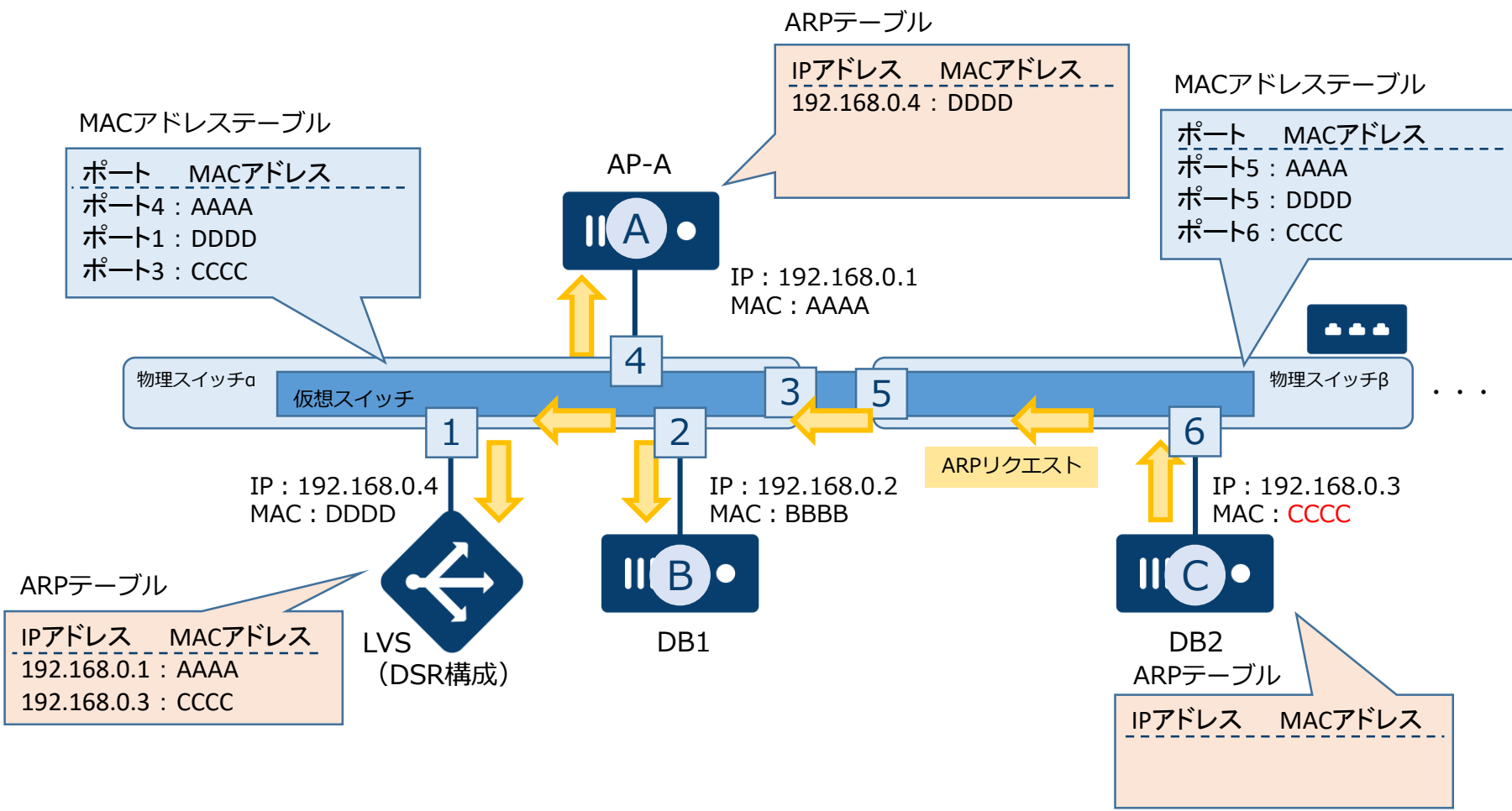
LVSはサーバCのMACアドレスを知ることができたので、サーバCに通信することができます。





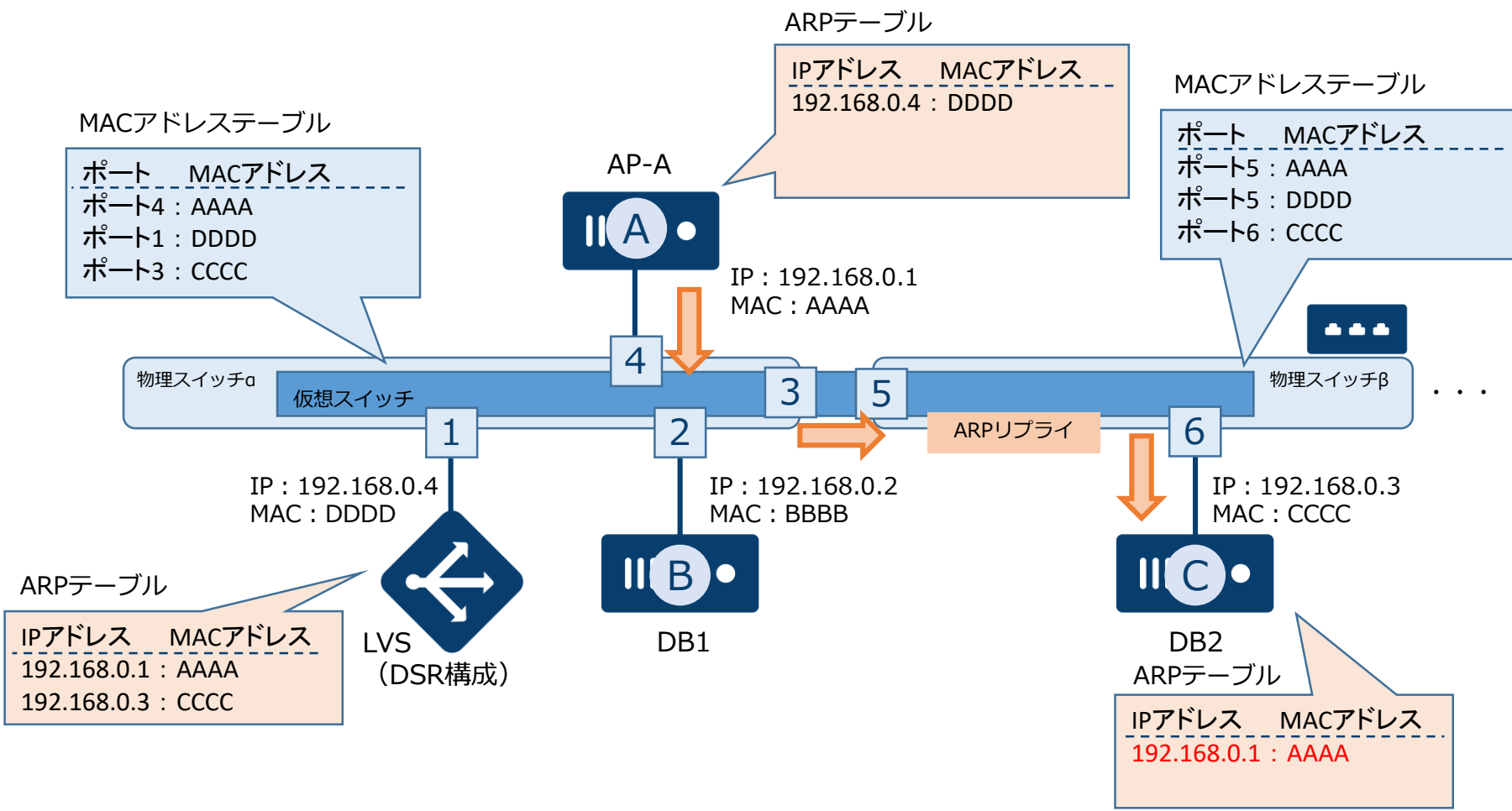
DSR構成のため、サーバCからの戻りの通信は、LVSを経由せず直接サーバAに送信されます。

サーバCはサーバAのMACアドレスを調べるため、ARPリクエストを送信します。



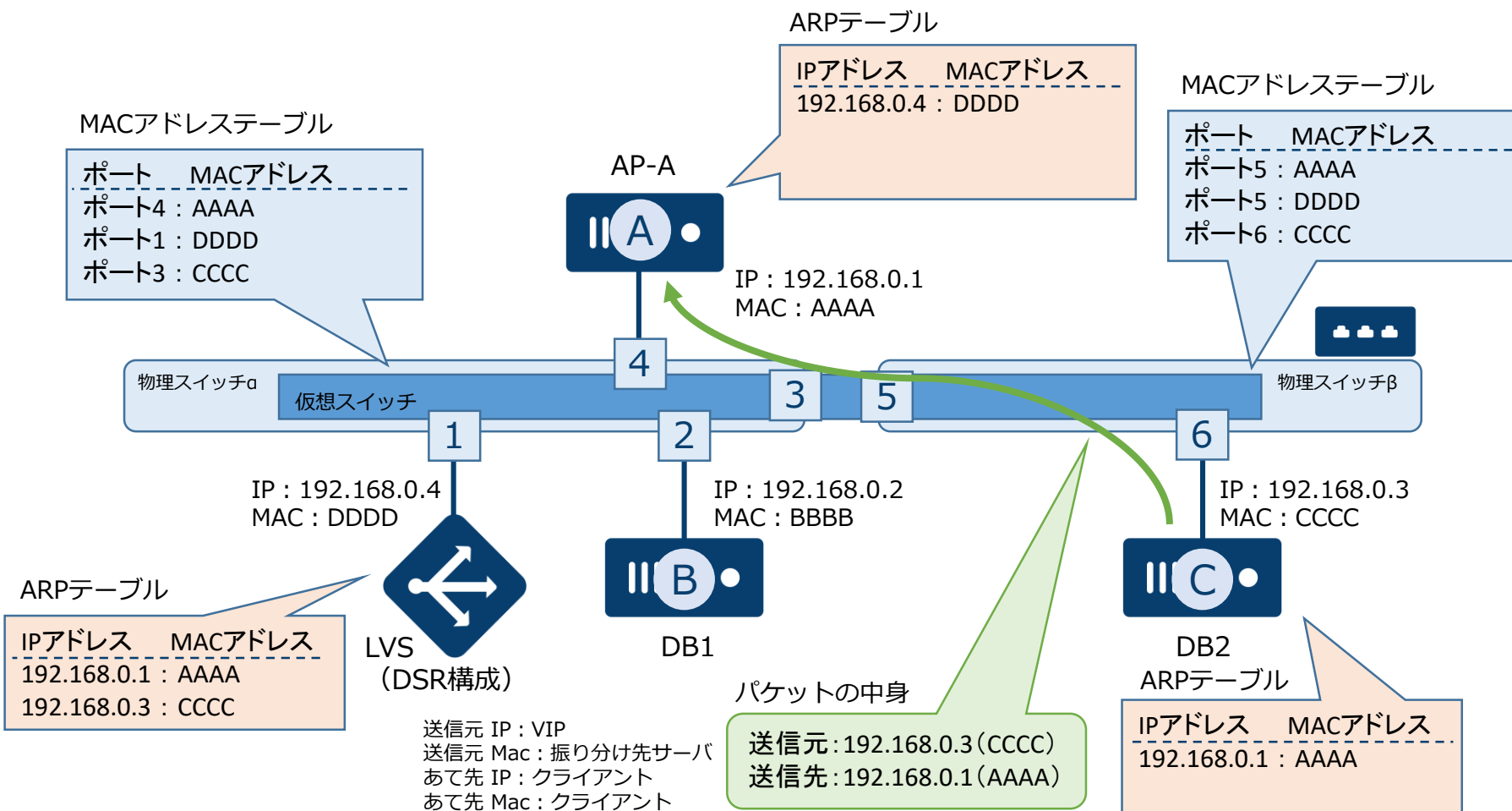


ARPリプライがサーバCに届き、サーバCはサーバAのMACアドレスを知ることができました。

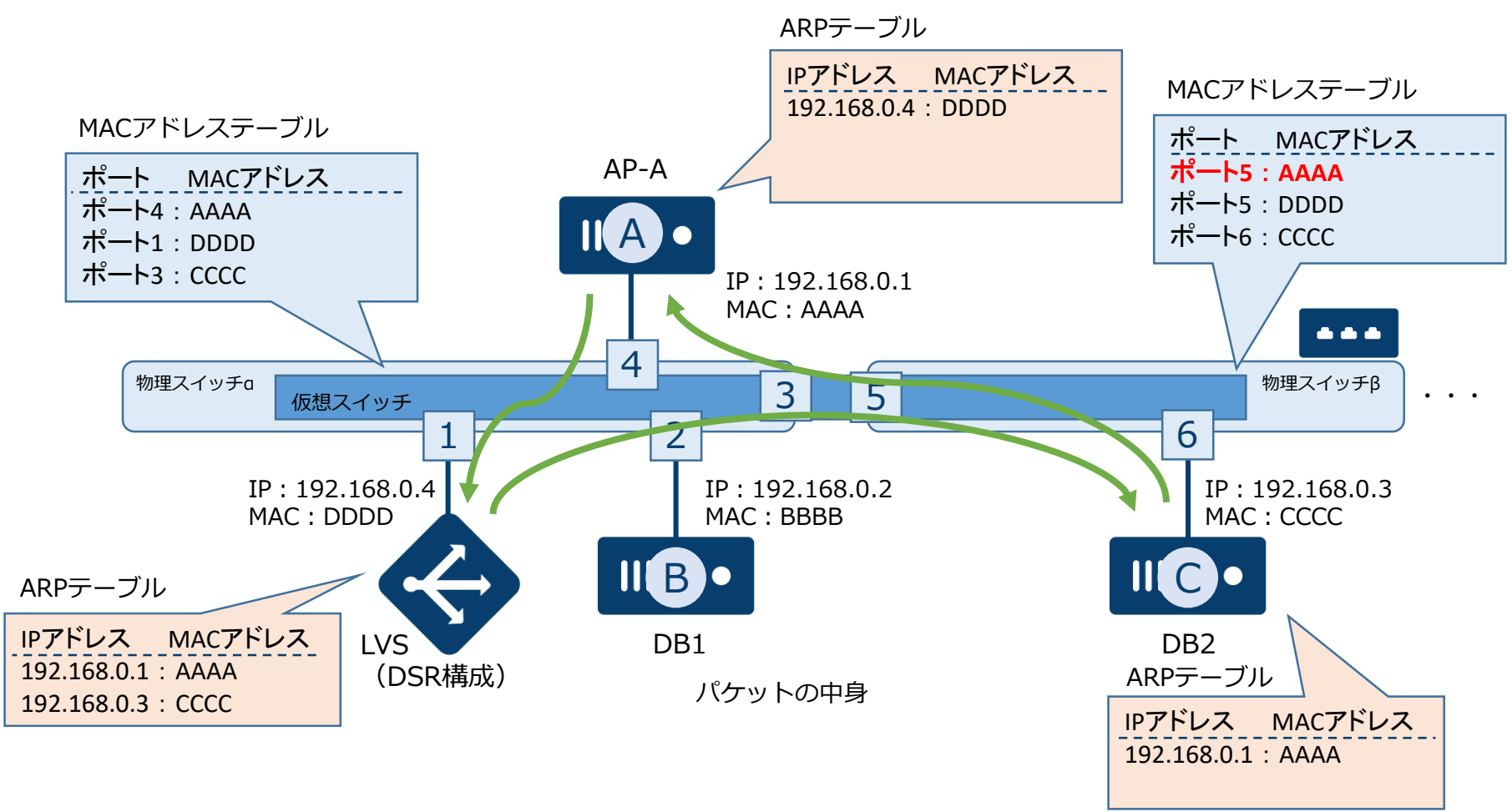




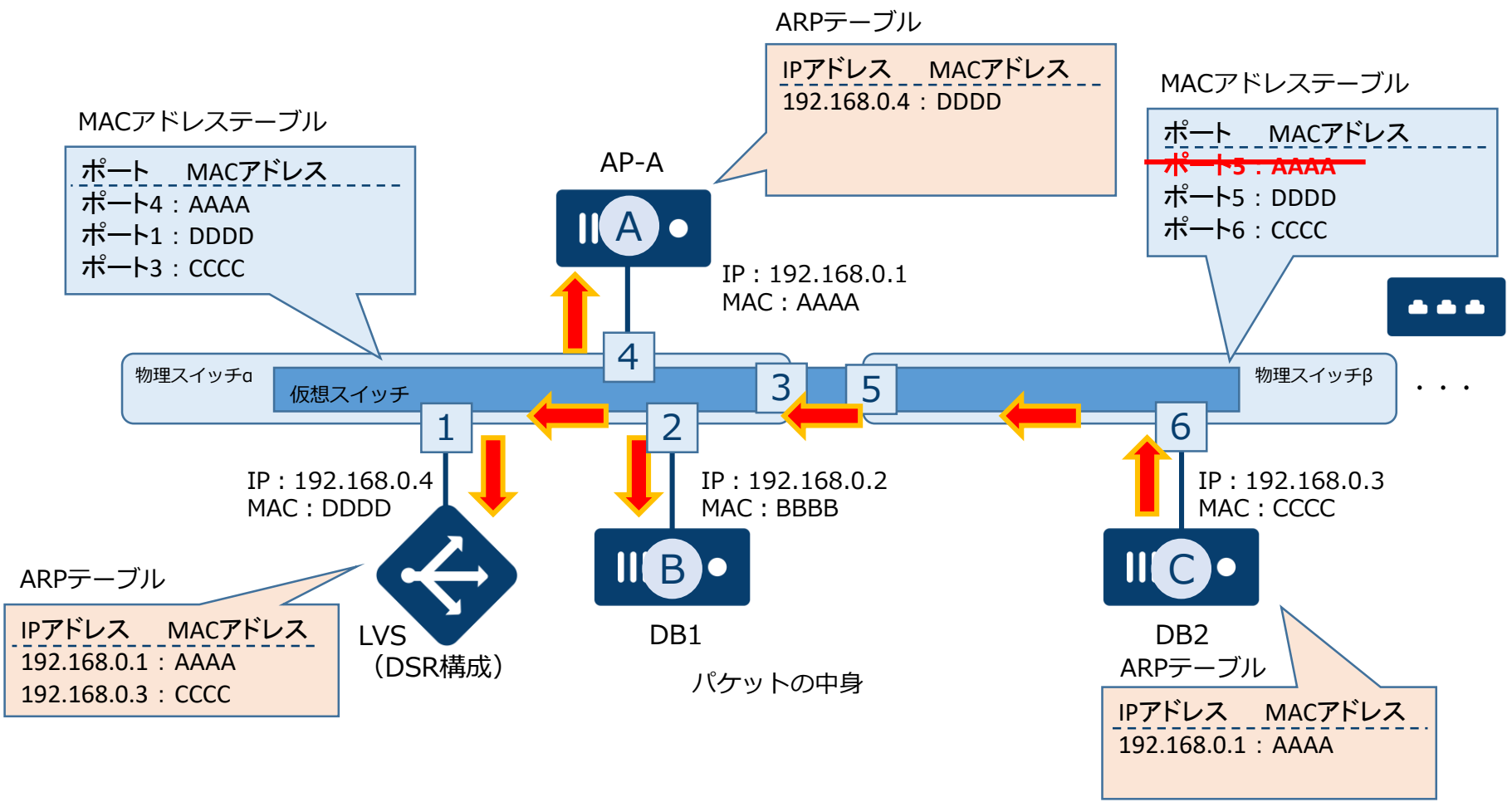
サーバCはサーバAのMACアドレスを知ることができたので、サーバAに通信することができます。



このDSR構成ですとサーバAとサーバCによるARPリクエスト（ブロードキャスト）以外、「サーバA⇒サーバC」の通信が発生しません。また、LinuxサーバのARPテーブルは仕様上保持される傾向があり、サーバからARPリクエストは送信されません。つまり、物理スイッチβはサーバAのMACアドレスを、改めて知る機会をほぼ失ったこととなります。



この状態が続き、物理スイッチβのMACアドレステーブルの情報が消去（エージング）されると、サーバCからサーバA宛ての通信は、以下のようにフラッディングされます。





■ 対策例

各サーバから定期的なpingを打ち相互にMACアドレスの学習を促す仕組みの実装
(弊社のロードバランサー仮想アプライアンスは元々対策を実装済みです)