

並列分散コンピューティング (9)論理時計

大瀧保広

今日の内容

- 物理時計と論理時計
 - 事象の順序のみに着目した「時計」
- 事象と順序関係
- 論理時計
 - Lamportのタイムスタンプ
 - 論理ベクトル時間

物理時計と論理時計

- 通常「時計」というと、なんらかの物理量に基づく時計をイメージする。これを**物理時計**と呼ぶ。
(Clock-on-the-Wall などと表現されることもある)
- 論理時計**とは、システム内でだけ通用する、物理量に基づかない時計。
(すなわち物理世界の時間の流れと独立でかまわない)
- 分散システムの動作を考える際には、必要な**順序関係が正しく判断できれば十分**であることが多い。
- 順序関係を判断するだけであれば、**単調増加するだけのカウンタで十分**といえる。

$$t' > t \Rightarrow C(t') > C(t)$$

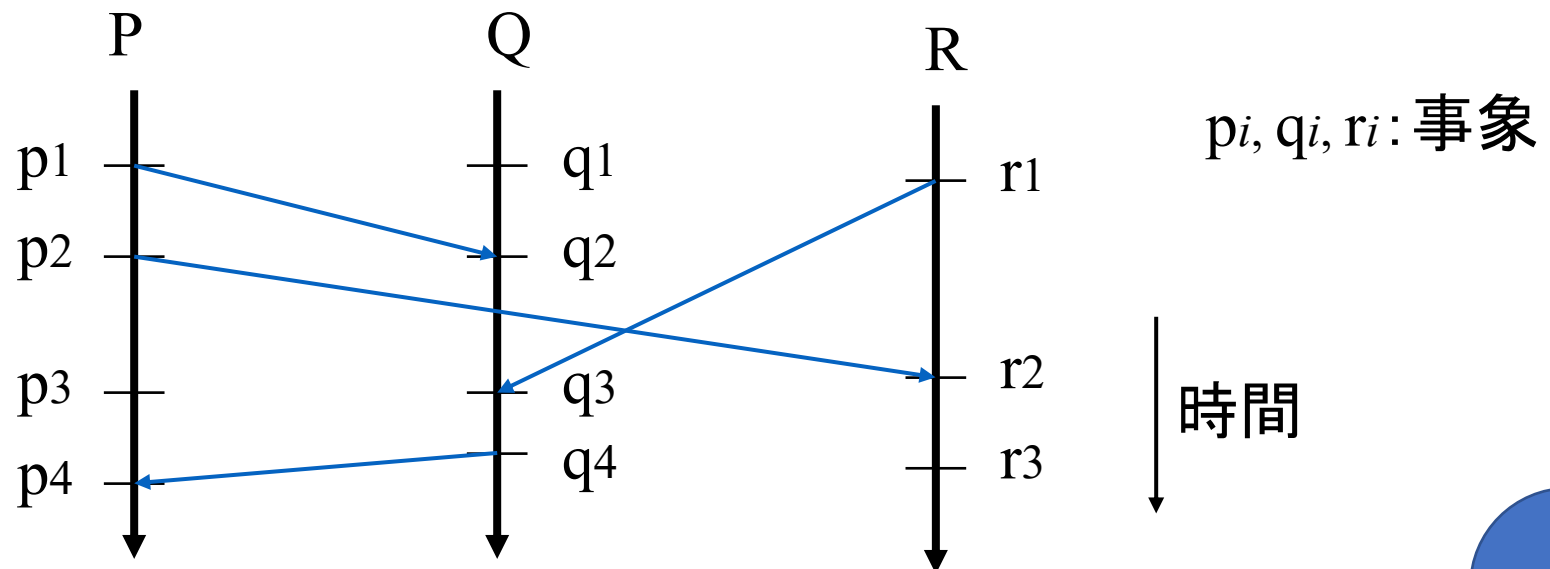
事象(イベント)と順序

事象の順序

- あるプロセス内で生じる様々な出来事を、**事象 (event)** と呼ぶ。
- 事象の性質
 - あるプロセスの事象が原因で、他のプロセスの事象が影響を受ける場合がある。
- 事象についての**仮定**
 - 同一のプロセス内のすべての事象は、発生した順に**順序付けることができる**。（発生順を認識できる）
 - 事象は**瞬時に終了**する。（事象は時間軸上の点となる）
 - 同一のプロセス内では**同時に2つの事象は起きない**。
- 分散システムにおいて重要なことは、複数プロセスにおける事象の順序について正しく**合意**できることである。

時空ダイアグラム(Time-Space-Diagram)

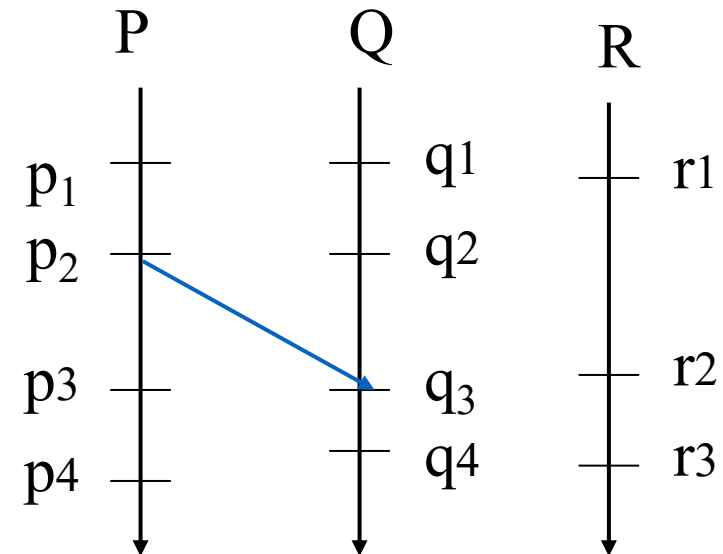
- プロセス間でのメッセージのやり取りを表現する方法
 - 水平方向にプロセスを配置し、垂直方向に時間を表す。
- 事象(event)の起きた時点に横棒（とラベル）をつける。
- プロセス間の矢印は、メッセージの送受信を表す。



事象の順序関係

事象の順序関係を記号 \rightarrow で表す。

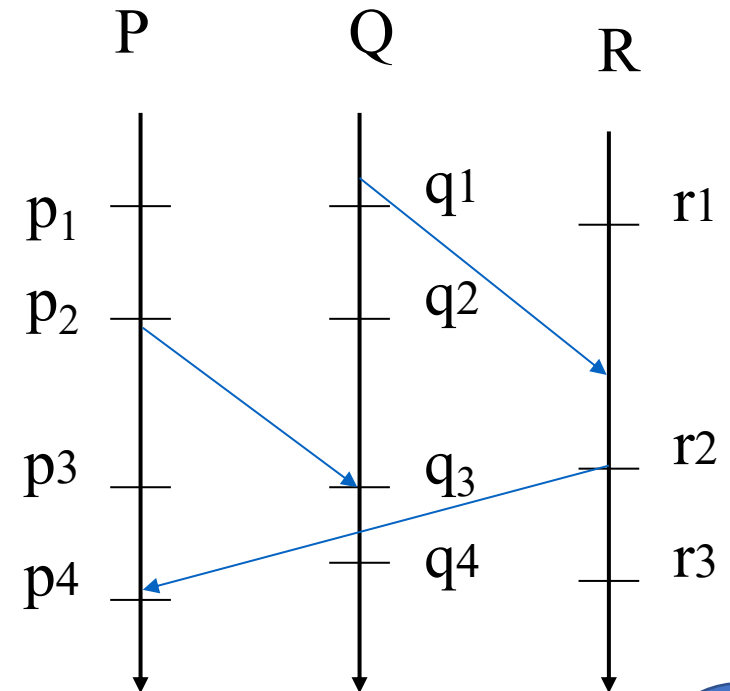
- 事象aが事象bより先に起こるとき、 $a \rightarrow b$ と書く。
- 同一のプロセス内では事象の順序関係は明らか。
例： $p1 \rightarrow p2$
- あるメッセージについて、aが送信したという事象でbがそれを受信したという事象ならば、 $a \rightarrow b$ である。
例： $p2 \rightarrow q3$
- 順序関係には推移律が成り立つ。すなわち、事象a, b, cについて $a \rightarrow b$ かつ $b \rightarrow c$ ならば $a \rightarrow c$ である。



事象の順序関係 (つづき)

- 2つの事象 a, b の間に、 $a \rightarrow b$ と $b \rightarrow a$ のいずれも成り立たないとき、
事象 a と b は並行 (concurrent) であるといい、 $a \parallel b$ と書く。

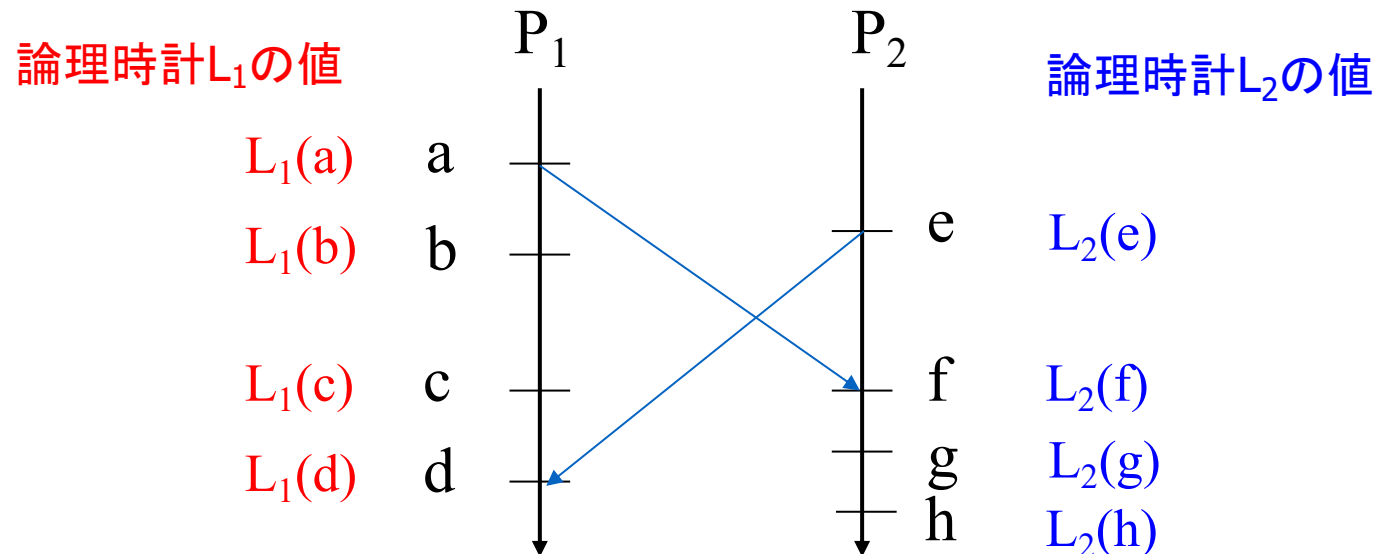
- 同一プロセス内の先行関係：
 $p_1 \rightarrow p_2, q_2 \rightarrow q_3, r_1 \rightarrow r_2$ など
- メッセージの送受信による関係：
 $p_2 \rightarrow q_3, q_1 \rightarrow r_2, r_3 \rightarrow p_4$ など
- 推移律によって導かれる関係：
 $p_1 \rightarrow q_3, q_1 \rightarrow p_4$ など
- 並行
 $q_2 \parallel r_1, p_3 \parallel q_4$ など



事象と論理時計

論理時計

- プロセス P_i における論理時計を論理時計 L_i とする。
(論理時計はプロセスごとにある)
- プロセス P_i における事象 a の論理時計 L_i の値を $L_i(a)$ と表す。

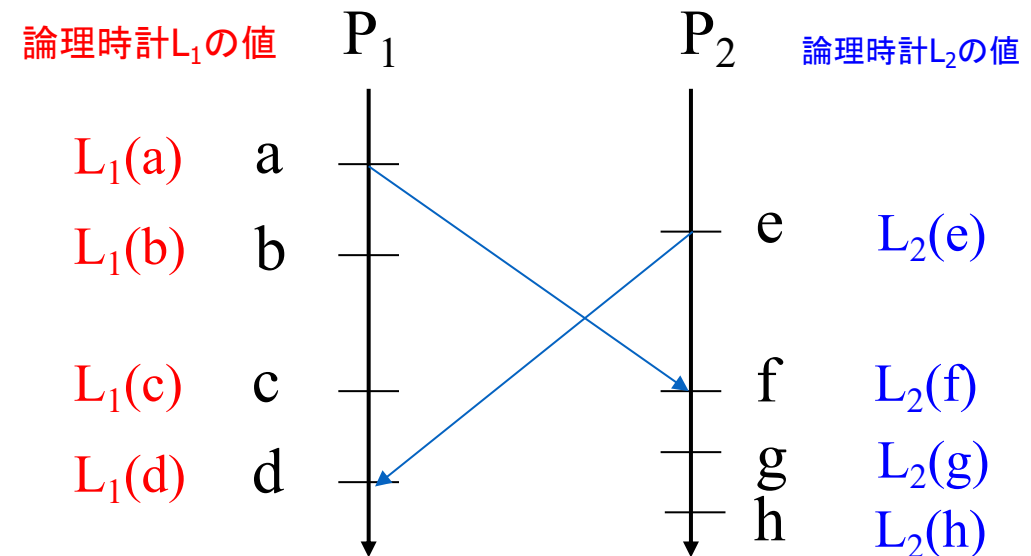


論理時計の要件ver1

事象の順序関係に基づき、
後に発生したイベントほど大きい値であること。

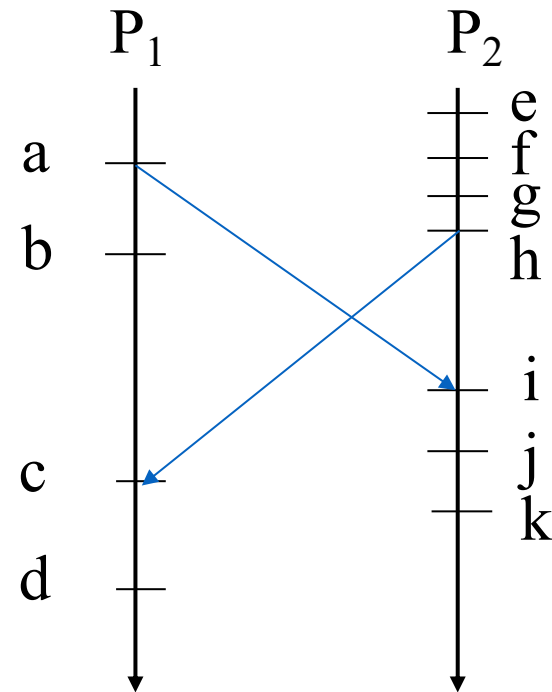
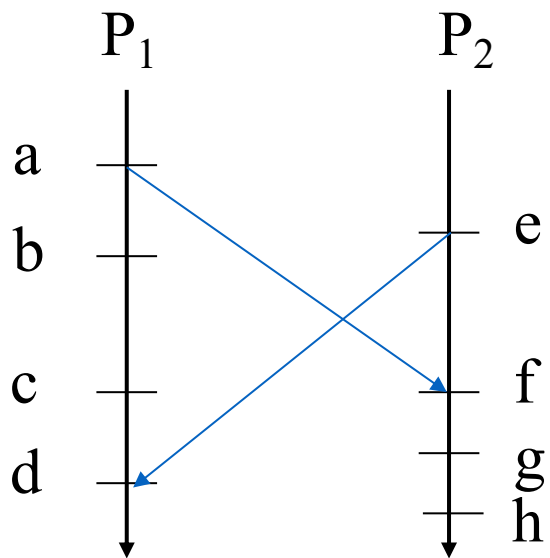
■あるプロセス P_i において、
 $a \rightarrow b$ ならば $L_i(a) < L_i(b)$

■ a がプロセス P_i における
メッセージ送信事象で、
 b がそのメッセージの
プロセス P_j における
受信事象ならば、
 $L_i(a) < L_j(b)$



要件を満たすような値 $L(\cdot)$ は
どうやったら実現できるだろうか？

単純に番号を振ってみる？



Lamportのタイムスタンプ

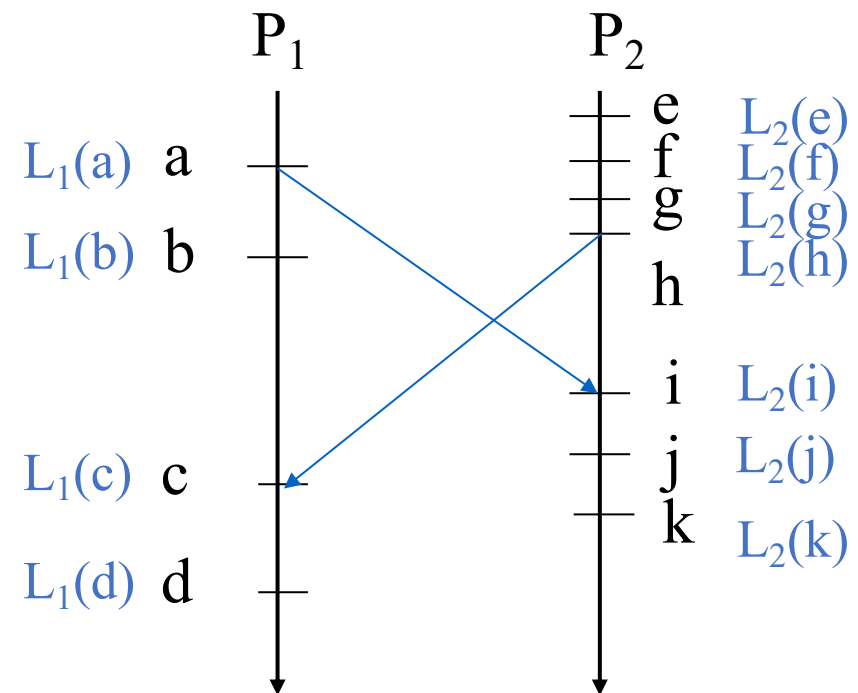
基本的な考え方

■ L_i : プロセス P_i の論理時計

■ $L_i(e)$: 事象 e のタイムスタンプ
(=時刻)

■ $L_1(a) < L_1(b)$ を実現するには？

論理時計 L_1 の値を増やし、
その値をタイムスタンプとして
付与すればよい。



基本的な考え方（つづき）

■ $L_2(h) < L_1(c)$ を実現するには？

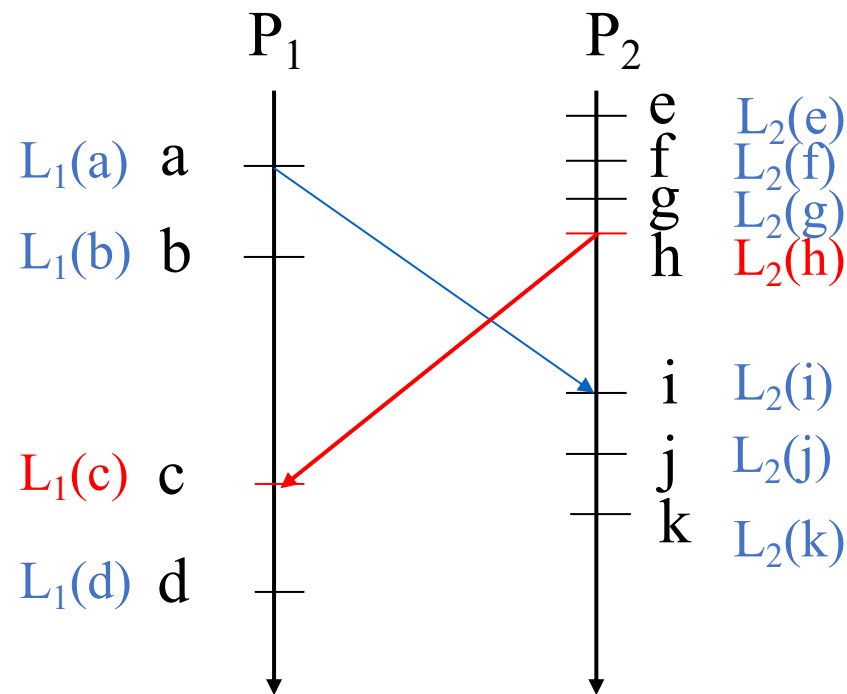
■ 正確には

$L_2(h) < L_1(c)$ と $L_1(b) < L_1(c)$
を実現するには？

■ 論理時計 L_1 の値を

$L_1(b)$ と $L_2(h)$ のどちらよりも
大きい値まで増やし、
その値をタイムスタンプとして
付与すればよい。

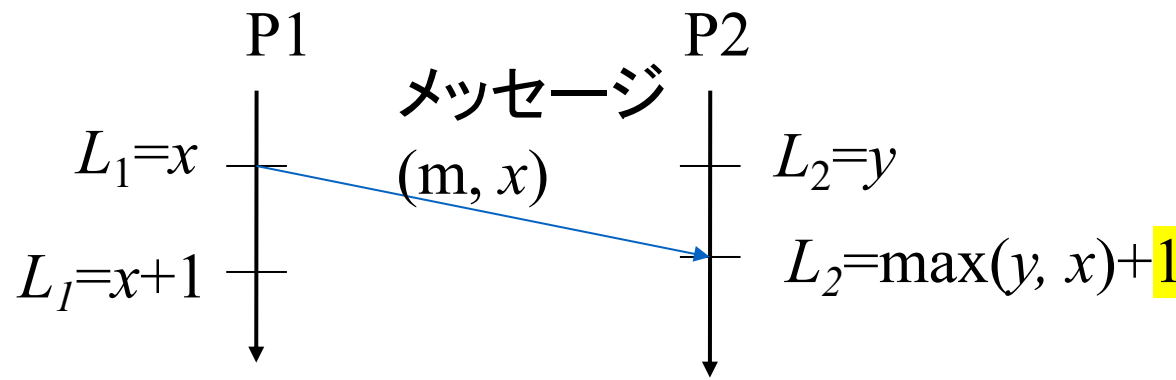
■ しかし、 P_1 では $L_1(b)$ はわかるが、
 $L_2(h)$ の値が不明である。
→メッセージにタイムスタンプを
つけて教えてもらう。



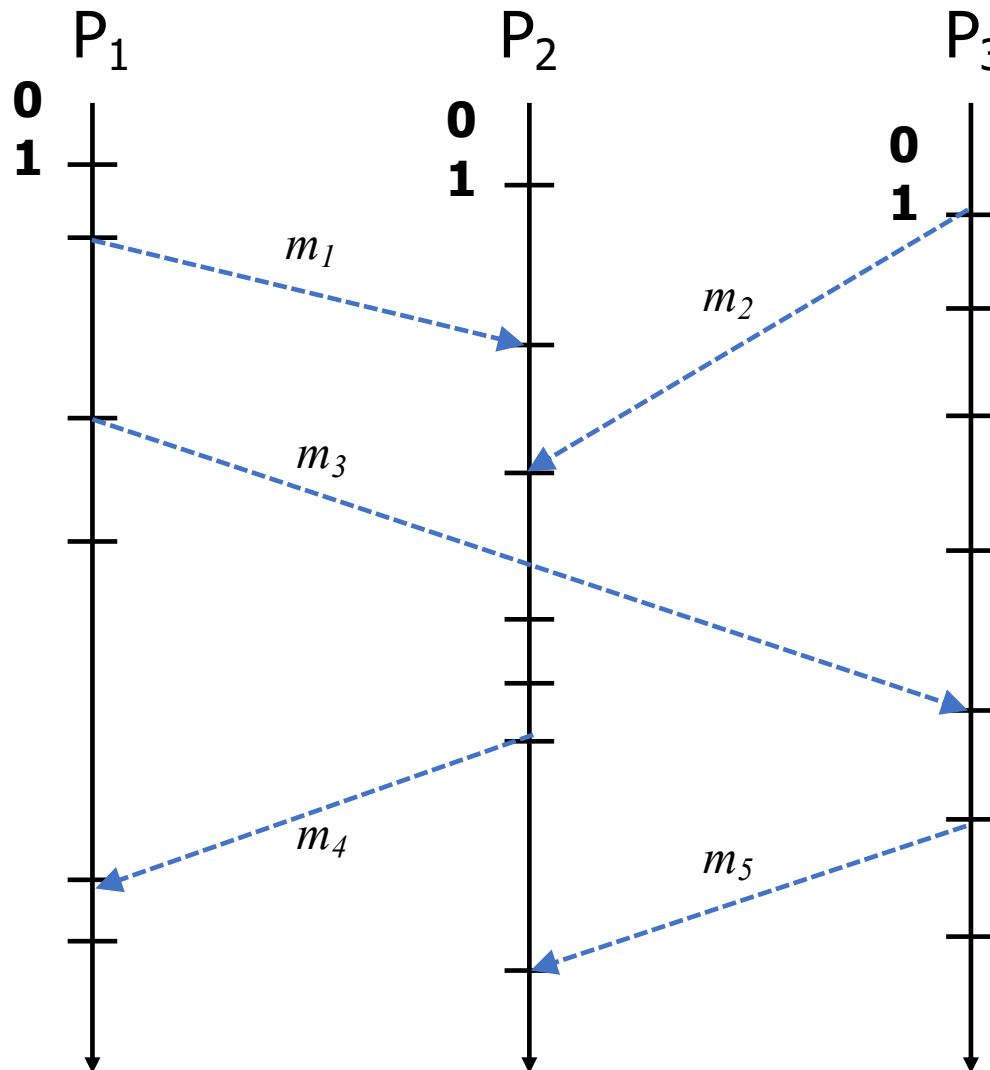
Lamportのタイムスタンプ

- 各プロセス P_i で各自の論理時計 L_i を0に初期化
- 各プロセス P_i では:
事象が発生するたびに $L_i = L_i + 1$; $L_i(\text{event}) = L_i$
- プロセス P_i がプロセス P_j にメッセージ m を送信するとき
 - 送信側 P_i :
 $L_i = L_i + 1$; $t = L_i$, message = (m, t) ; $L_i(\text{send event}) = L_i$
 - 受信側 P_j :
 $L_j = \max(L_j, t)$; $L_j = L_j + 1$; $L_j(\text{receive event}) = L_j$

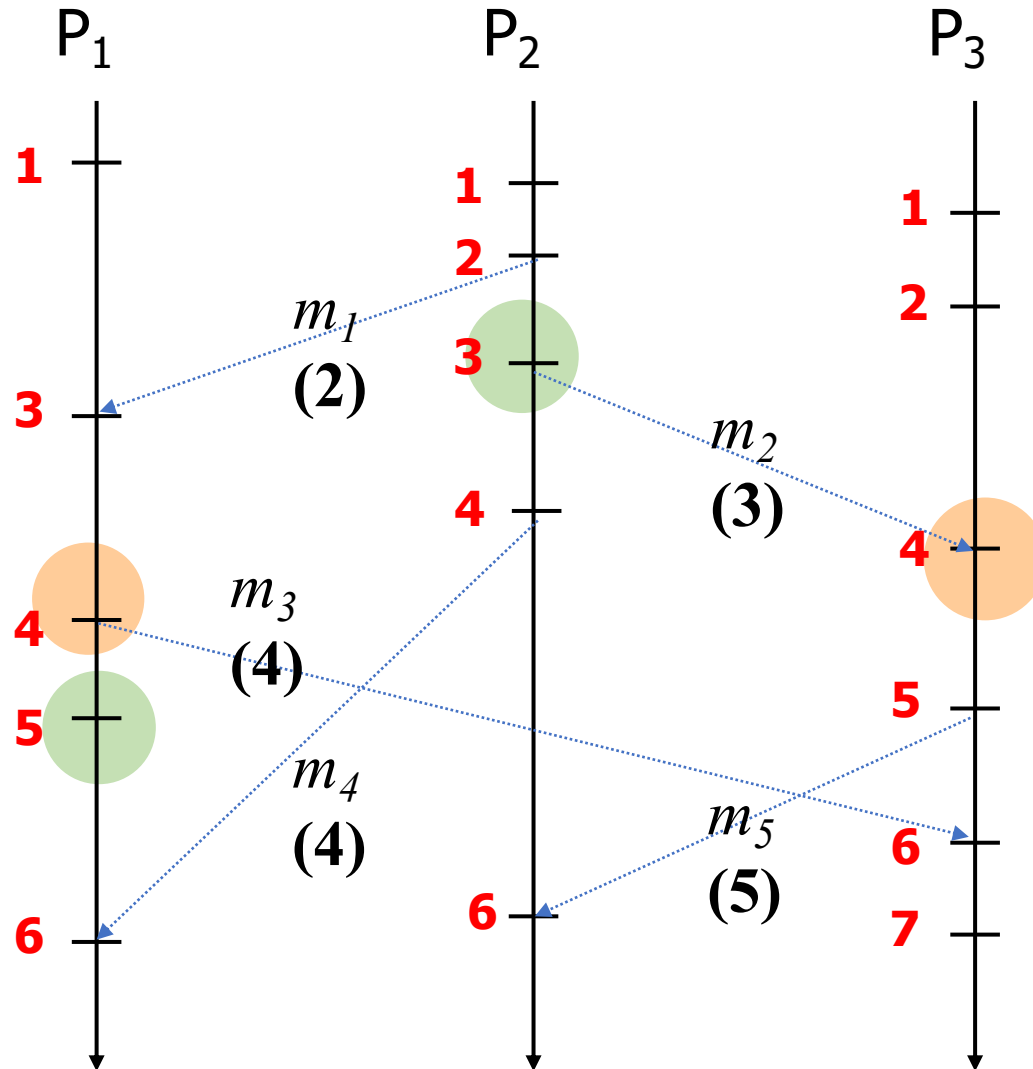
正の数であれば
何でもよい



演習：Lamportのタイムスタンプ



Lamportのタイムスタンプの問題点



Lamportのタイムスタンプの問題点

「 $a \rightarrow b$ ならば $L(a) < L(b)$ である」ようには構成することができた。しかし...

**タイムスタンプの値の大小関係から
事象の順序関係が判断できない。**

■問題点1：

$L(a) < L(b)$ だからといって、 $a \rightarrow b$ とは限らない。

■問題点2：

$L_i(a) = L_j(b)$ のとき、事象 a, b の順序はどう考えればいい？
（「同時に発生した」という意味？）

Lamportのタイムスタンプの問題点

- タイムスタンプの値が事象の順序を意味しない。

$L_1(a) < L_2(b)$ だからといって、 $a \rightarrow b$ とは限らない。

- これは、実際には順序関係のないところ、すなわち $a \parallel b$ であるところ に $L_1(a) < L_2(b)$ といった順序関係をつけてしまったということ。
→ 「並行に処理可能だった」ことがわからなくなっている。

- 分散システムでは、論理時計の値に基づいて順序関係を判断したいのに、これでは困る。

論理時計の要件ver2

■ 次のような要件を満たせる論理時計が欲しい

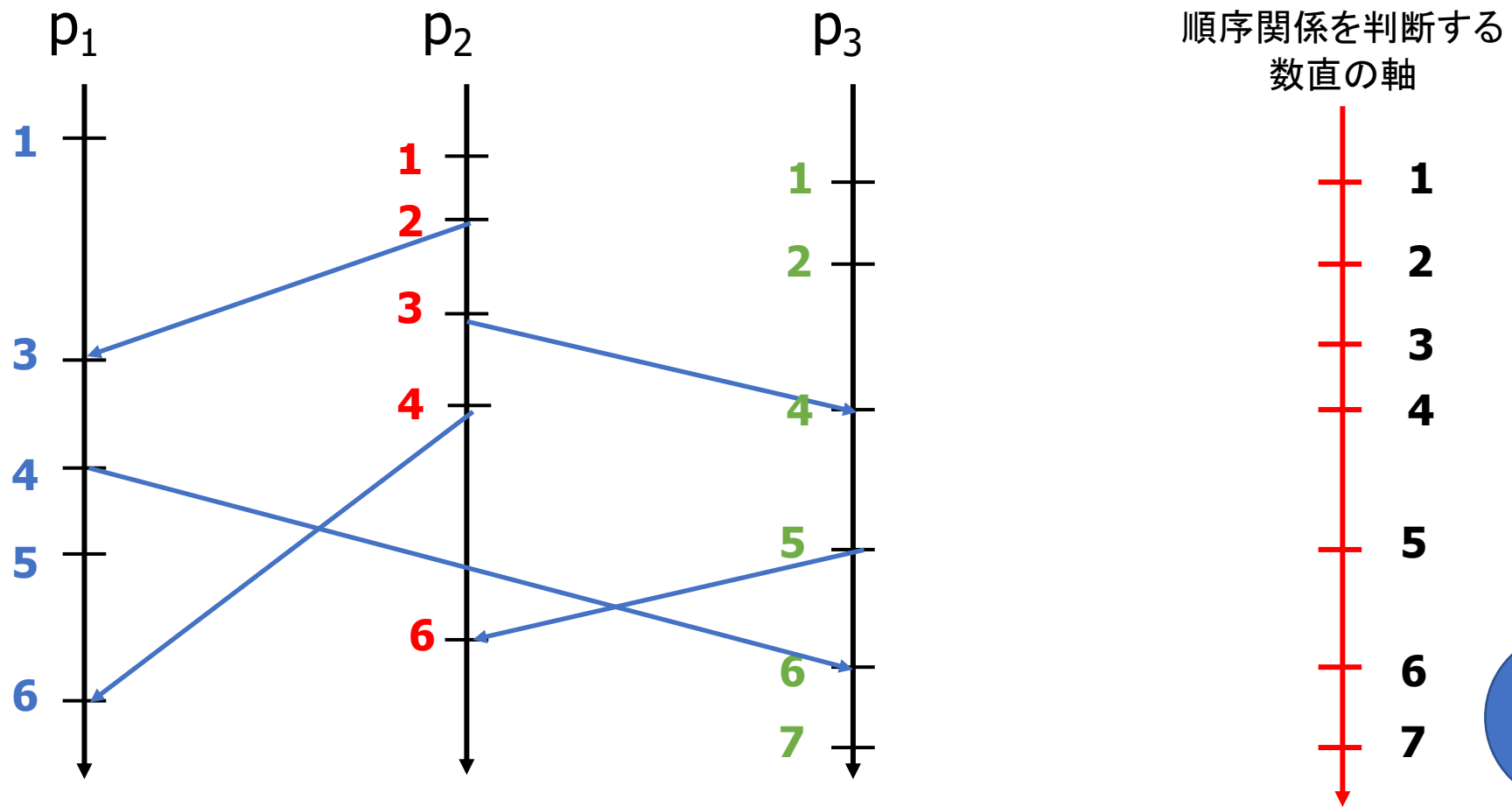
1. $a \rightarrow b$ ならば $L(a) < L(b)$ であること。
2. 「 $L(a) < L(b)$ ならば $a \rightarrow b$ 」といえること。
3. $L(a)$ と $L(b)$ をみて、 $a || b$ であるかどうか判断できること。

■ どうやったら構成できるのだろうか？

■ Lamportのタイムスタンプで要件2 と 要件3が実現できなかった真の原因はどこなのか？

Lamportのタイムスタンプの問題の原因

- 各プロセスではそれぞれ固有の論理時計でカウントしているのに、区別せずに扱っていること。

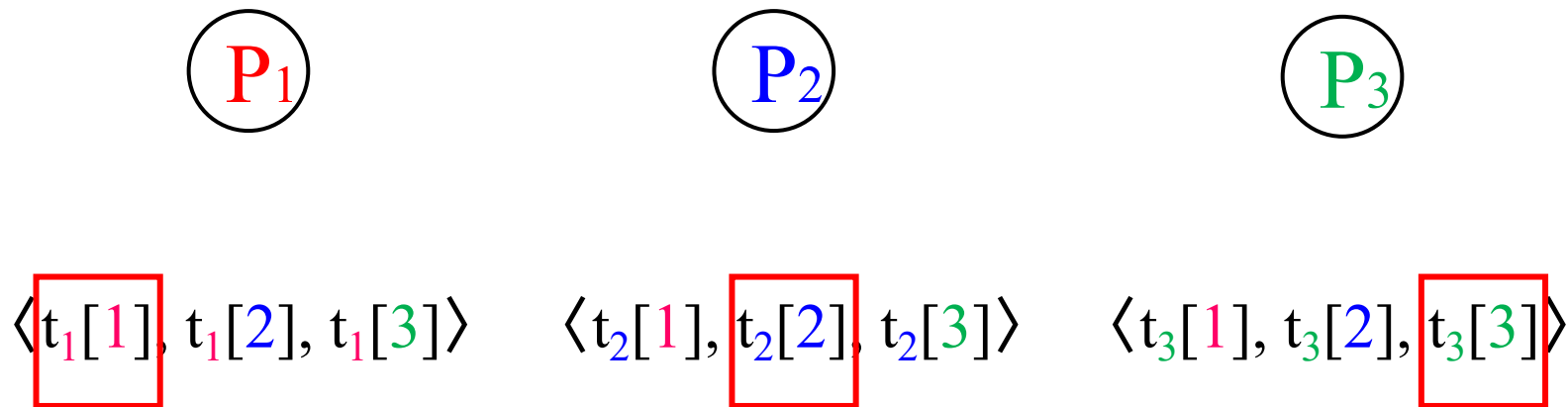


論理ベクトル時間

1つの数値（スカラー）で表すのではなく、
数値の組み（ベクトル）で表す

論理ベクトル時間

- 各プロセスの論理時計の値を一緒にせず、常に区別する。
- 各プロセスの値を組にして、一つのベクトルとして扱う。
- システムとしては、このベクトルをタイムスタンプとして使用する。（＝論理ベクトル時間）



$t_i[j]$: プロセス P_i が保持するプロセス P_j の論理時計の値

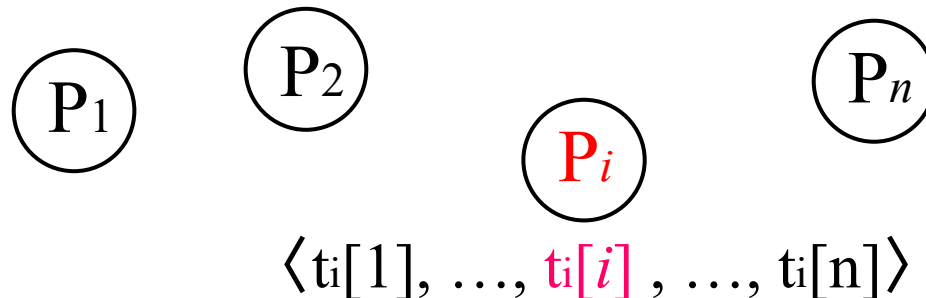
論理ベクトル時間

n プロセスからなるシステムの場合、
論理ベクトル時間は **n 個の数値の組み** のベクトル で表す。

プロセス P_i がもつ論理ベクトル時間
 $t_i = \langle t_i[1], t_i[2], \dots, t_i[n] \rangle$

$t_i[j]$: プロセス P_i が保持する プロセス P_j の論理時計の値

$t_i[i]$: P_i 自身の論理時計の値



論理ベクトル時間の処理(1)

[メッセージの送受信]に関係しない事象の場合

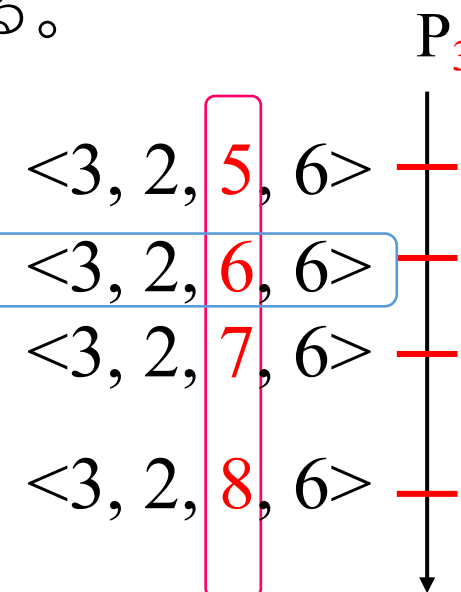
■ プロセス P_i で事象が起きたら：

1. P_i は自分自身の論理時計を進める。

$$t_i[i] \leftarrow t_i[i] + 1$$

2. 事象に論理ベクトル時間を付与する。

イベントに付与される
論理ベクトル時間
(タイムスタンプ)



P3の論理時計

重要： P_3 が保持する他のプロセスの
論理時計の値は変化しない。

論理ベクトル時間の処理(2)

[メッセージの送受信]に関わるイベントの場合

- 送信イベント：通常どおり論理ベクトル時間を付与し、メッセージに「送信イベントの論理ベクトル時間」をつけて送信。
- 受信イベント：受信プロセスでは論理ベクトル時間を次のように更新し、受信イベントに付与する。

受信プロセス P_i の論理ベクトル時間の各要素 $t_i[h]$ について

$$t_i[h] = \begin{cases} \max(t_i[h], t_j[h]) & (h \neq i) \\ t_i[h] + 1 & (h = i) \end{cases}$$

ここで $t_j[h]$ は受信した「送信イベントの論理ベクトル時間」の要素。

論理ベクトル時間の処理(2)

$$t_i[h] = \begin{cases} \max(t_i[h], t_j[h]) & (h \neq i) \\ t_i[h] + 1 & (h = i) \end{cases}$$

$\langle 4, 2, 3, 4 \rangle$
 $\langle 5, 2, 3, 4 \rangle$
 $\langle 6, 4, 6, 4 \rangle$

P1の論理時計

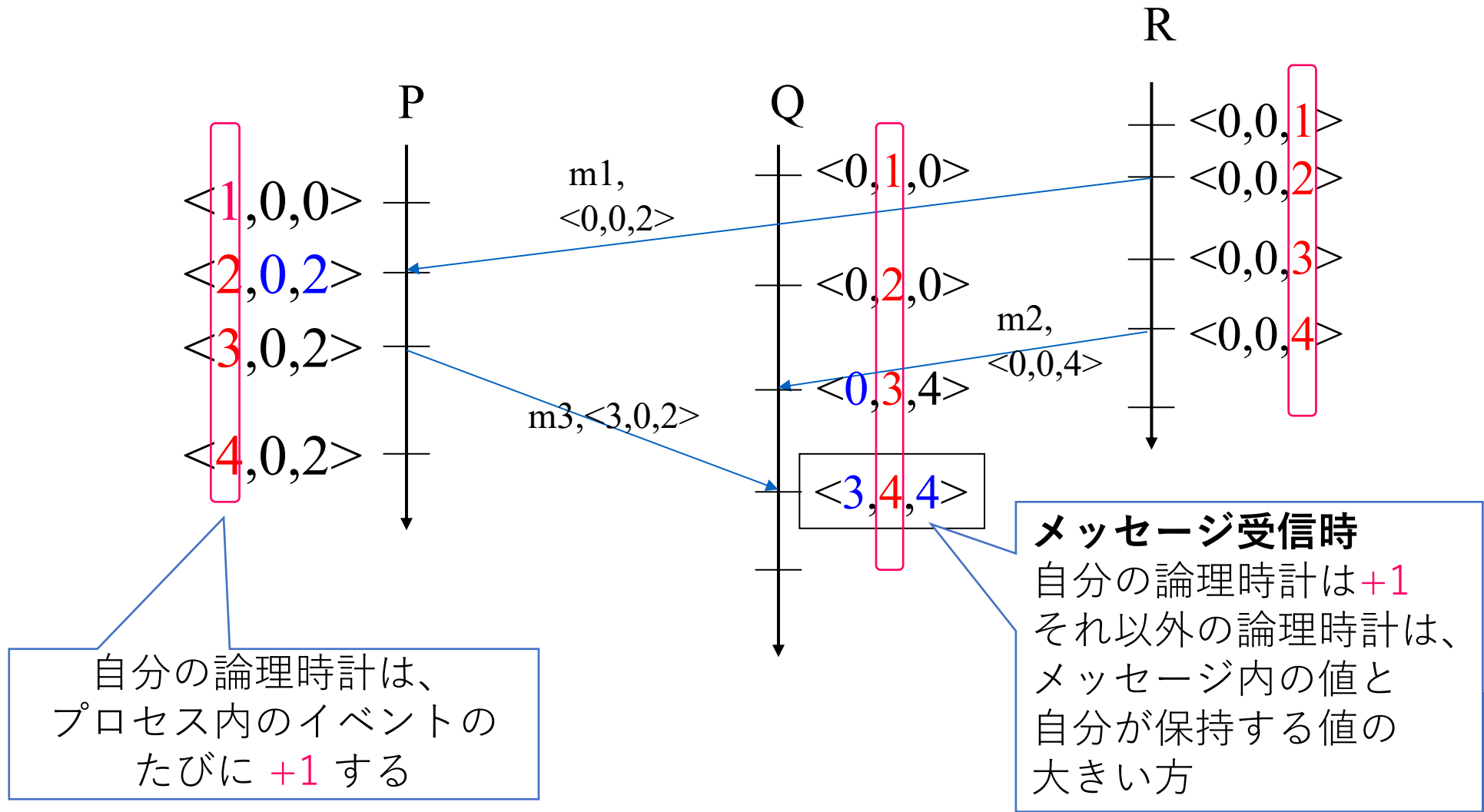
P₁

m, $\langle 3, 4, 6, 1 \rangle$

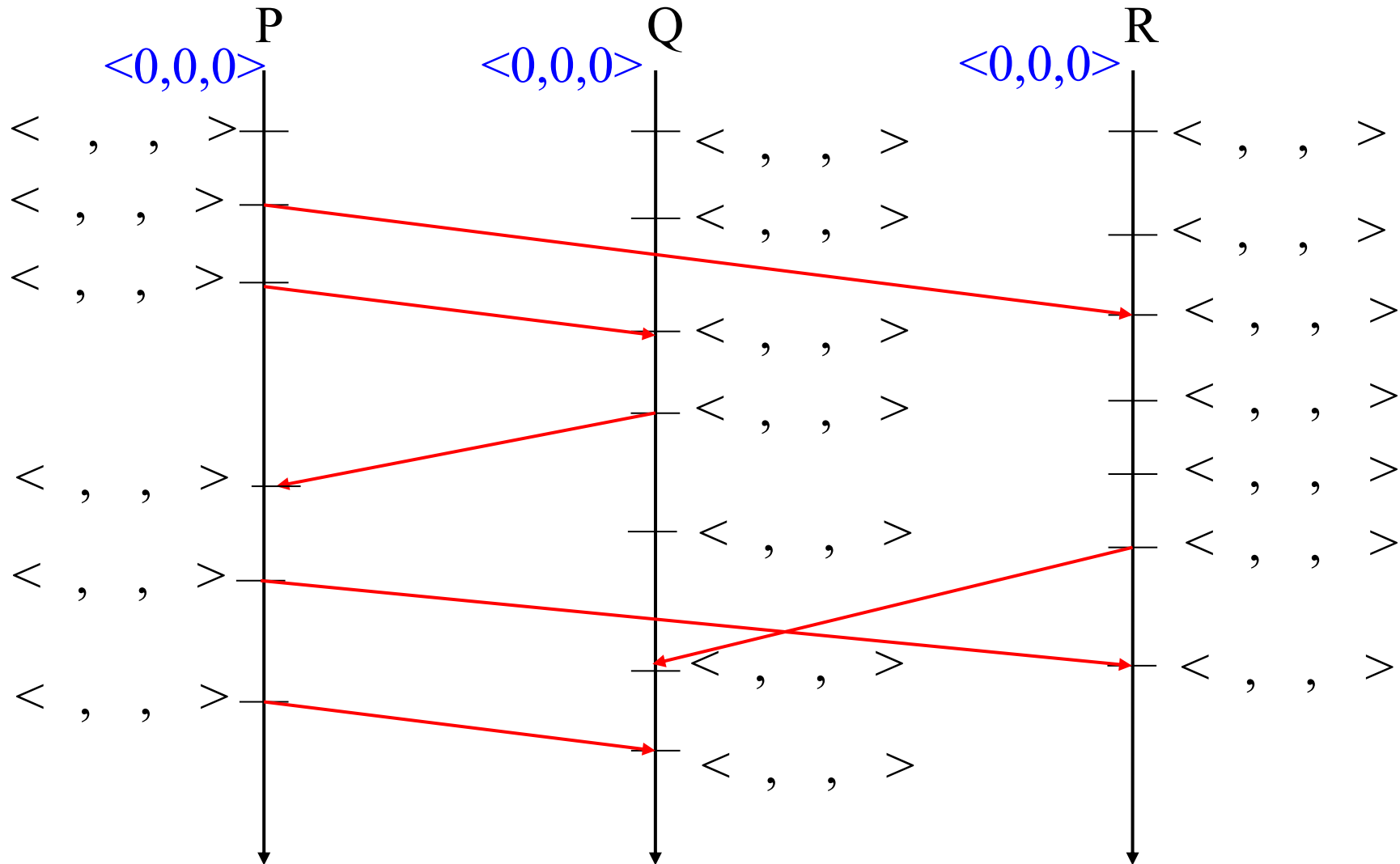
P₃

$\langle 3, 4, 5, 1 \rangle$
 $\langle 3, 4, 6, 1 \rangle$

論理ベクトル時間の例



演習：論理ベクトル時間



論理ベクトル時間の大小関係

■ 2つの論理ベクトル時間

$V = (v_1, v_2, v_3, \dots, v_n)$ と $V' = (v'_1, v'_2, v'_3, \dots, v'_n)$ があるとき、

■ 「 $V < V'$ である」とは、

全ての $i \in \{1, 2, 3, \dots, n\}$ に対して $v_i \leq v'_i$ であり、かつ、

ある $i \in \{1, 2, 3, \dots, n\}$ に対して $v_i < v'_i$ が成り立つことである。

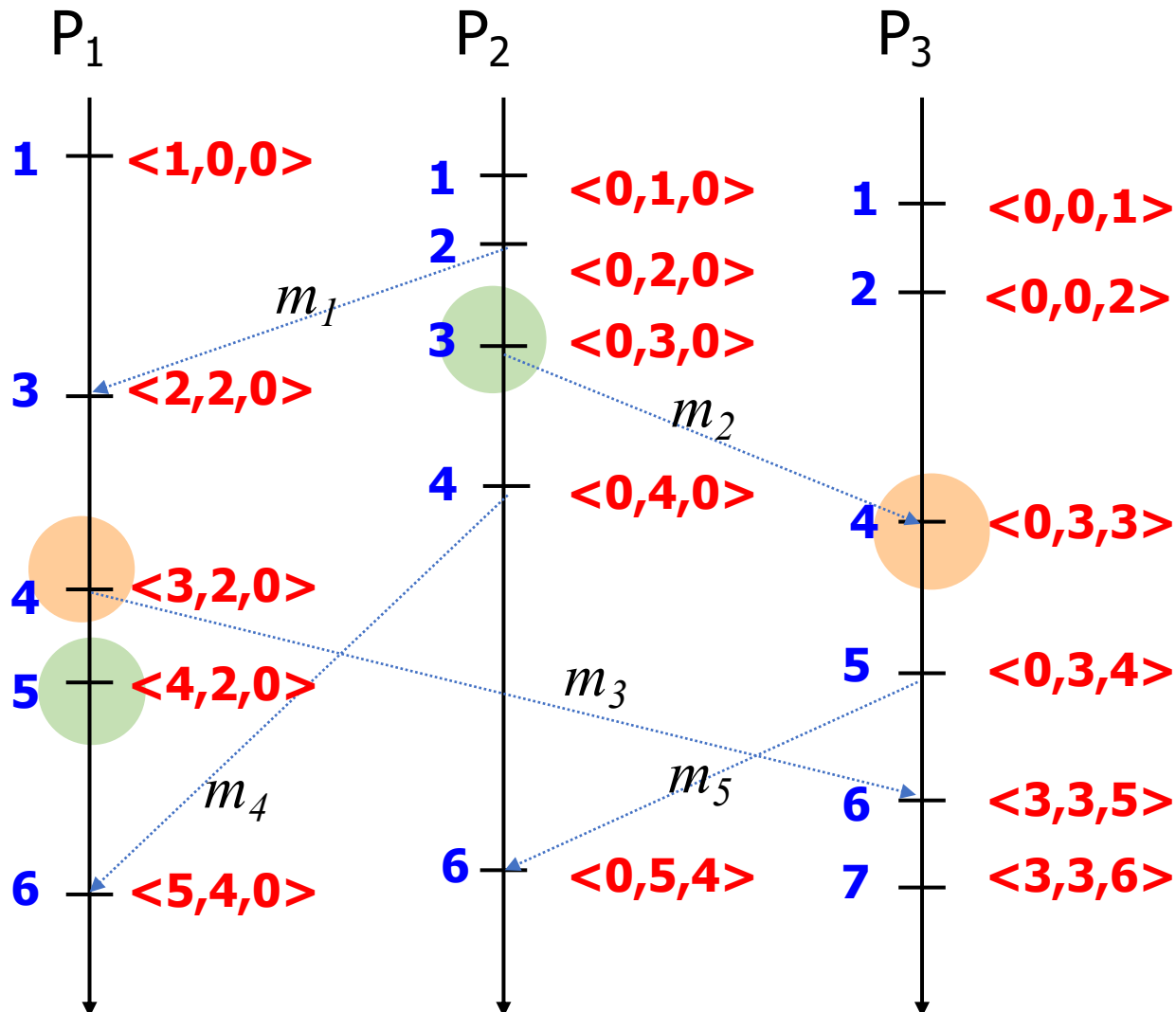
■ 成り立たない時には、 V と V' の間には大小関係は定義されない。

「 $\langle 3, 4, 5, 3, 2 \rangle$ 」 $<$ 「 $\langle 3, 5, 5, 3, 2 \rangle$ 」

「 $\langle 3, 5, 5, 3, 2 \rangle$ 」 $>$ 「 $\langle 3, 4, 5, 3, 1 \rangle$ 」

「 $\langle 3, 4, 5, 3, 2 \rangle$ 」 \parallel 「 $\langle 4, 4, 4, 3, 2 \rangle$ 」

LamportのTSで問題があった さっきの事例



$\langle 0, 3, 0 \rangle$ と $\langle 4, 2, 0 \rangle$
大小関係がない。
→ 並行

$\langle 0, 3, 3 \rangle$ と $\langle 3, 2, 0 \rangle$
大小関係がない。
→ 並行

過去イベント、未来イベント、並行イベント

あるイベント e に対して、それ以外の全てのイベントは、そのイベント e に対する
過去イベント、未来イベント、並行イベント
のいずれかに分類できる。

- 過去イベントの集合

$$\text{past}(e) = \{ e' \mid V(e') < V(e) \}$$

- 未来イベントの集合

$$\text{future}(e) = \{ e' \mid V(e') > V(e) \}$$

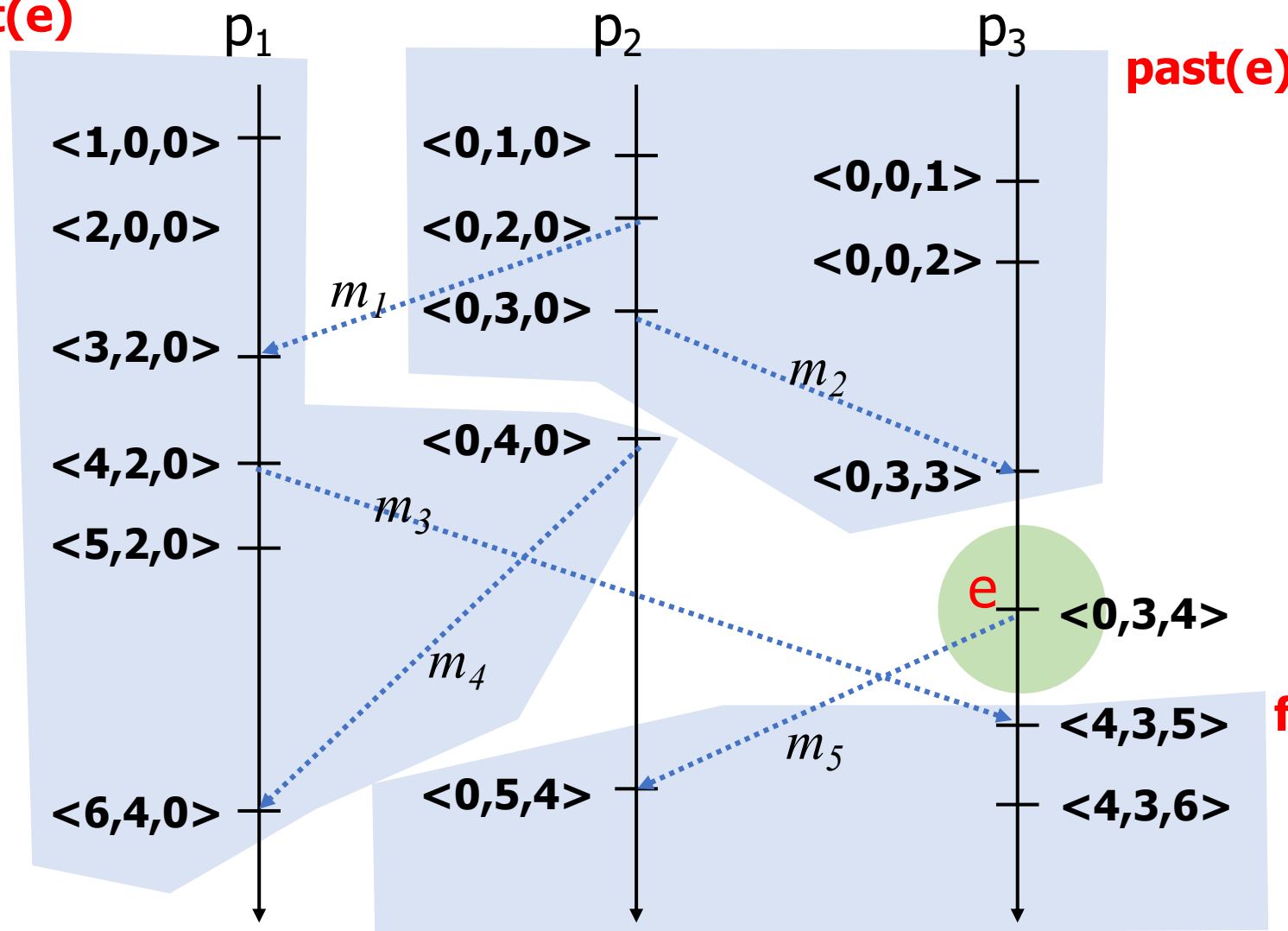
- 並行イベントの集合

$$\text{concurrent}(e) = \{ e' \mid (V(e') \not\leq V(e)) \wedge (V(e') \not\geq V(e)) \}$$

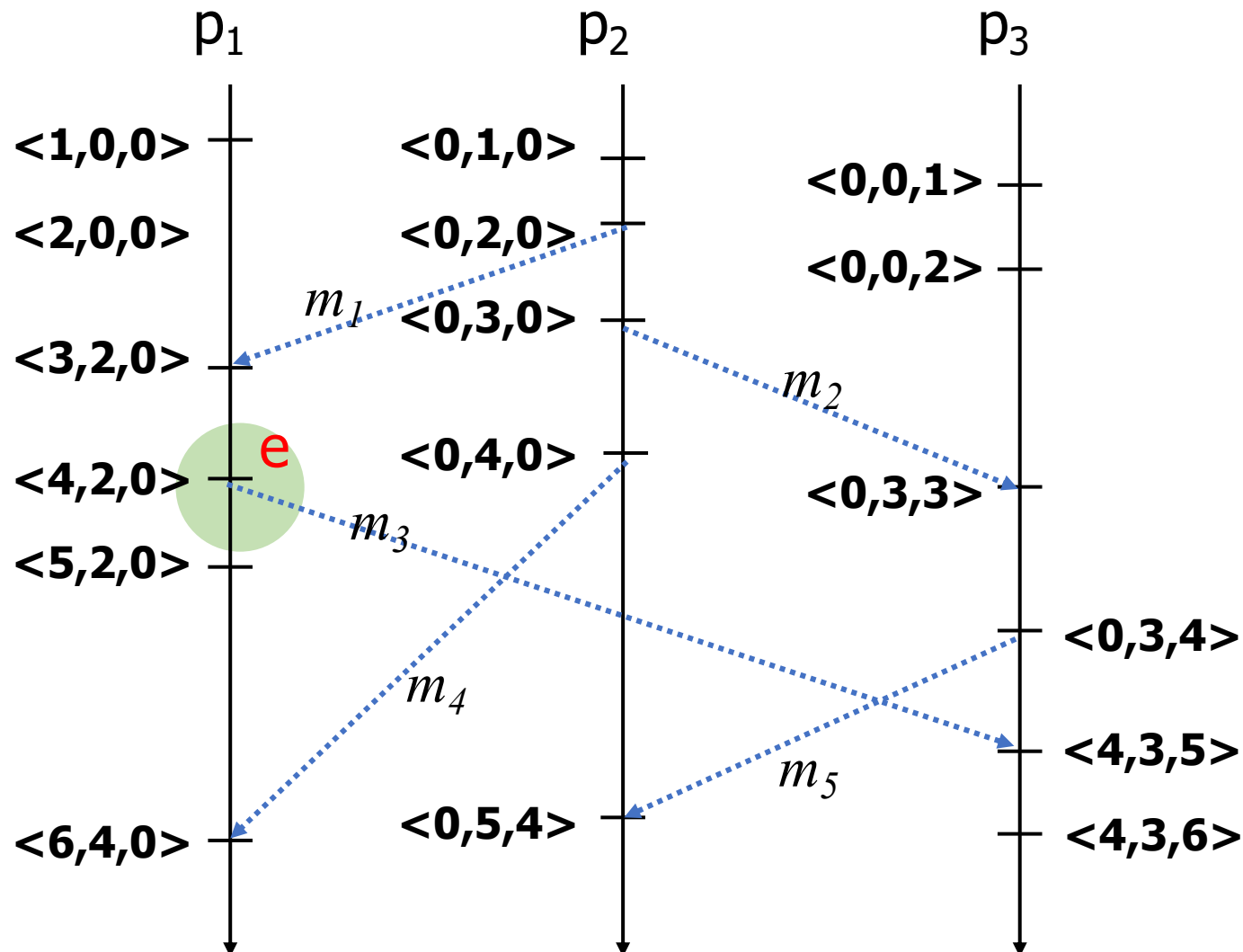
ここで $V(x)$ は事象 x の論理ベクトル時間。

過去イベント、未来イベント、並行イベント

concurrent(e)



演習：past(e), future(e), concurrent(e) の範囲を図示せよ



{過去 | 未来 | 並行}イベントは どう使う？

■原因の絞り込み

イベント e を実行中にエラーが見つかった場合、その原因は $\text{past}(e)$ の中にある。

■汚染範囲の特定

イベント e の時点でウィルスに感染したことが判明した場合、そのウィルスの影響を受けた可能性がある処理は、 $\text{future}(e)$ の範囲に限られる。
汚染範囲の特定

■スケジューリングの改善

$\text{concurrent}(e)$ のイベントは、 e と並行して処理することが可能なイベントである。

今日のまとめ

- 事象（イベント）の順序のみに着目
 - 前後と並行
- 論理時計
 - Lamportのタイムスタンプ
 - 問題点：タイムスタンプの値→順序関係 に難がある。
 - 論理ベクトル時間
- 論理ベクトル時間に基づき、あるイベントに対する過去イベント、未来イベント、並行イベントが決まる。
→分散システム運用上の判断基準