
めざせ！ Kafkaマスター
～Apache Kafkaで最高の性能を出すには～

2017年12月8日

株式会社日立製作所 OSSソリューションセンタ
伊藤雅博

- **伊藤 雅博 (いとう まさひろ)**

- **所属:** 株式会社日立製作所 OSSソリューションセンタ
- **業務:** Hadoop/Sparkを中心としたビッグデータ関連OSSの導入支援やテクニカルサポート

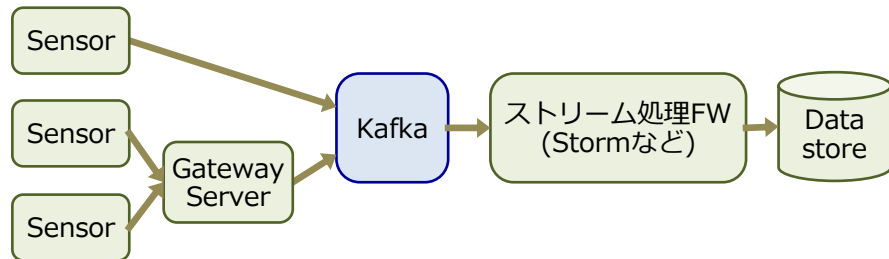
1. Apache Kafkaとは
2. Kafkaの内部構造とチューニングポイント
3. 性能検証の概要
4. 検証結果と考察
5. まとめ

1. Apache Kafkaとは

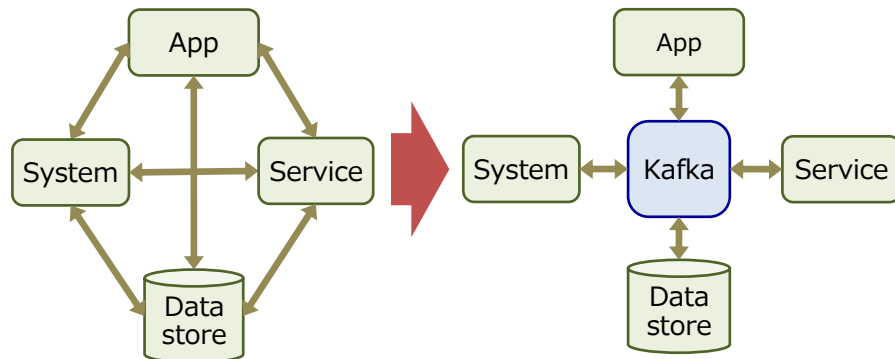
- スケーラビリティに優れた分散メッセージキュー
 - Pub/Subメッセージングモデルを採用
 - 書き込み/読み出し性能を重視しており、MQTTなどの標準プロトコルではなく、独自プロトコルを使用
- 主なユースケース:

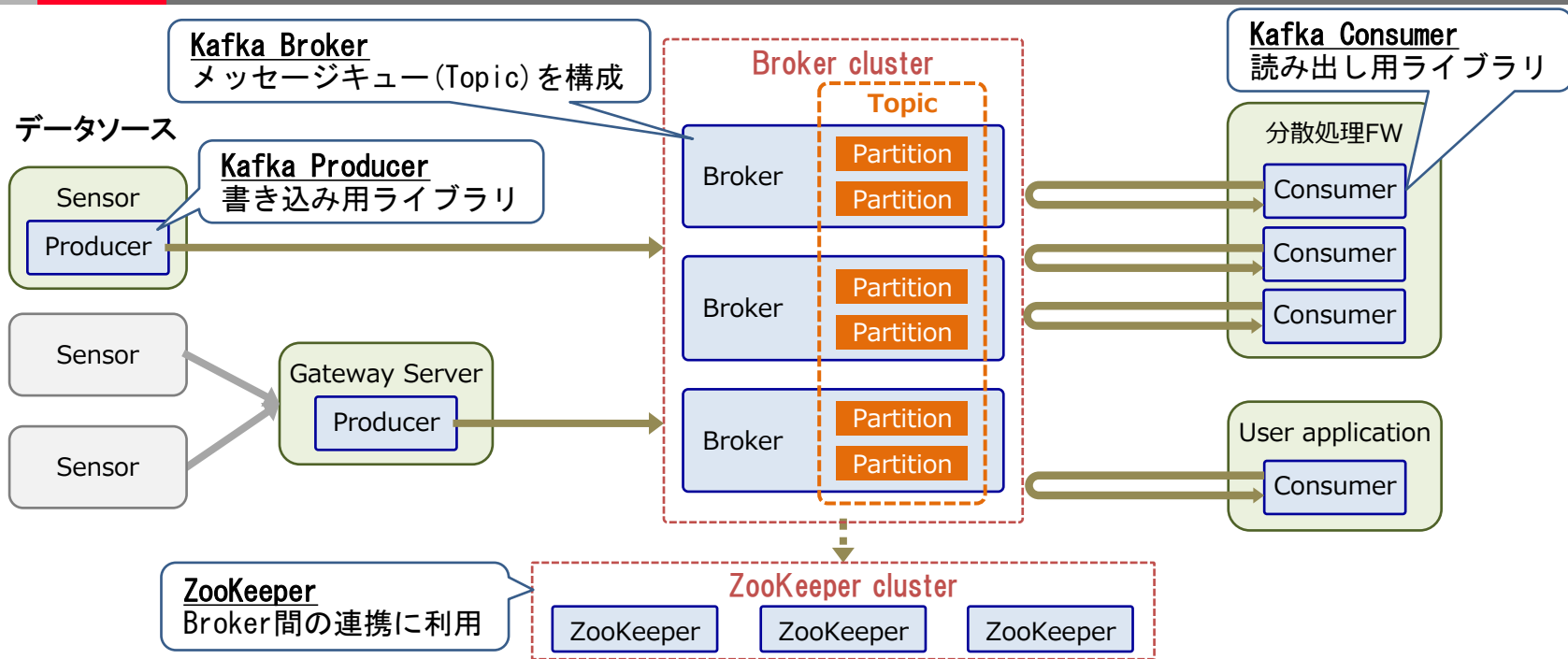
ストリームデータ処理システム におけるキューイング

データソース



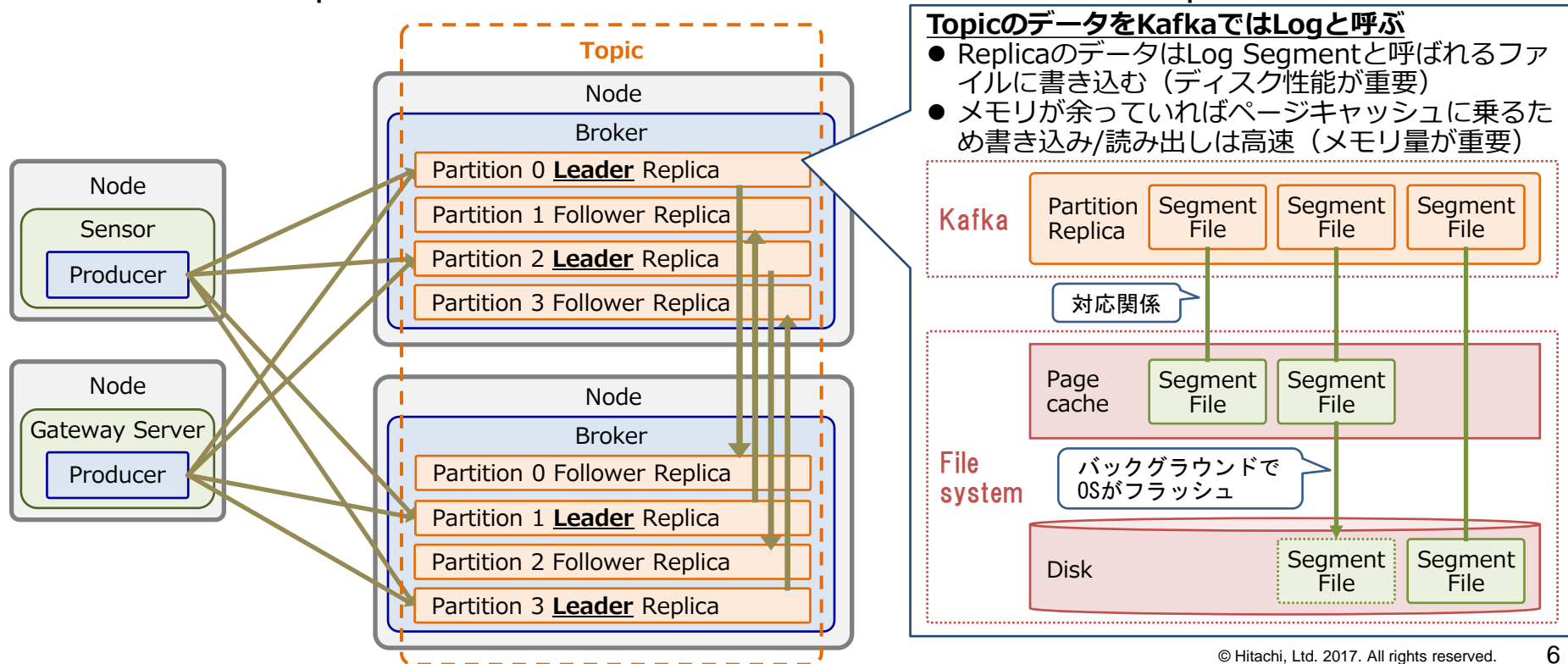
システム間で大量のデータを受け渡すためのパイプライン





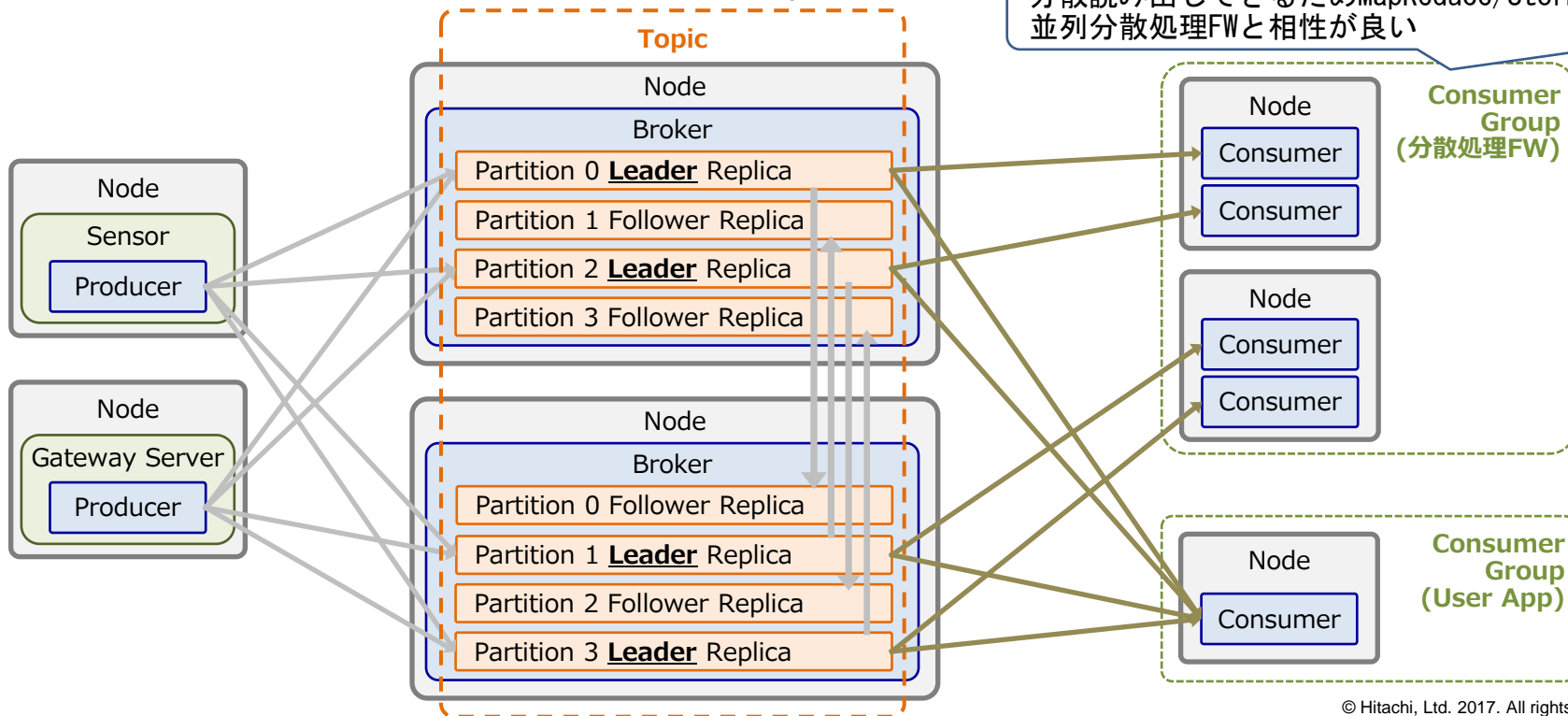
- Topic(1個の仮想的なキュー)を複数のBrokerに分散配置したPartitionで構成
- Topicのデータ書き込み/読み出しには Producer/Consumer ライブラリを使用
 - Java, C/C++, Pythonなど様々な言語用のライブラリが用意されている

- 各Partitionは 1個のLeader Replica + 0個以上のFollower Replica で構成
 - Leader Replicaにデータを書き込み、別Broker上の Follower Replica に複製する



- ConsumerはPartition単位でデータを並列読み出し可能
 - 複数のConsumerでグループを構成し、1Topicのデータをグループ内のConsumerで分散読み出し
 - 読み出し元は各 Partition の Leader Replica のみ

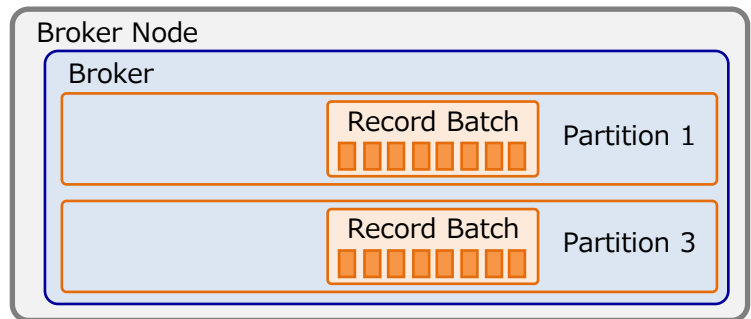
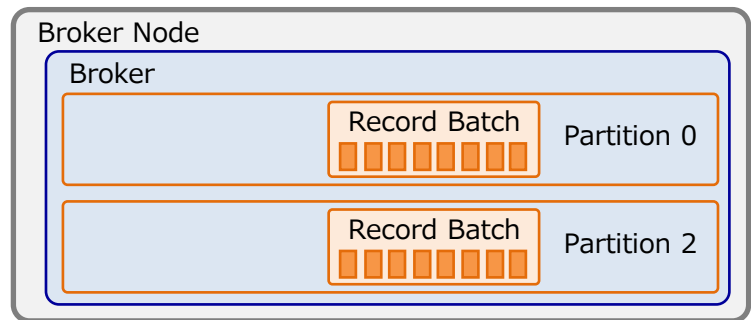
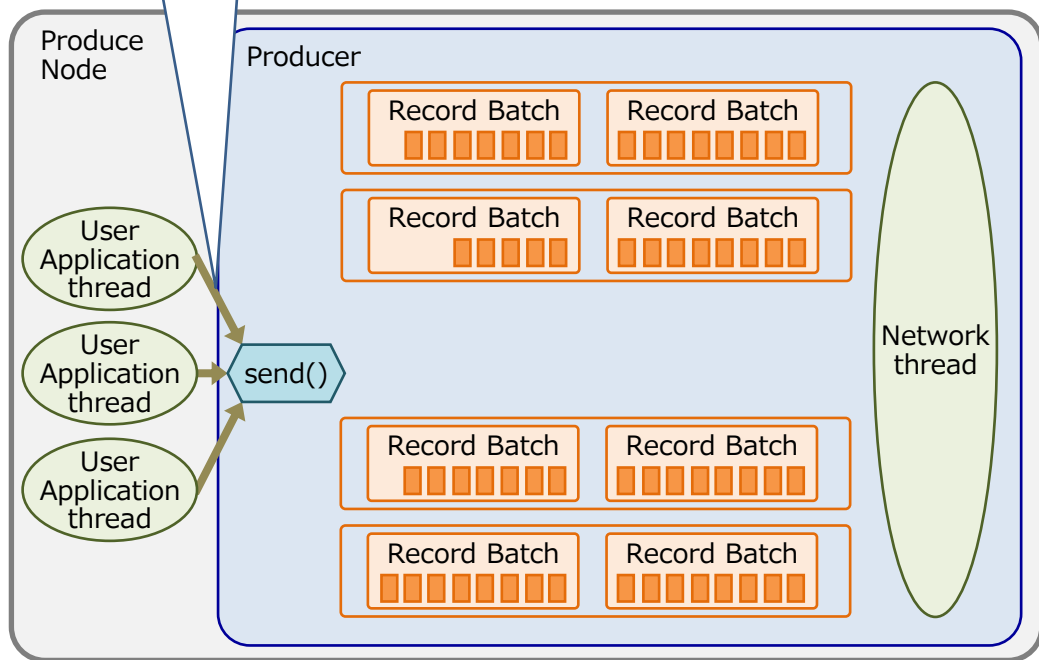
分散読み出しできるためMapReduce/Stormなどの
並列分散処理FWと相性が良い



2. Kafkaの内部構造とチューニングポイント

Producerの処理の流れ (1/5)

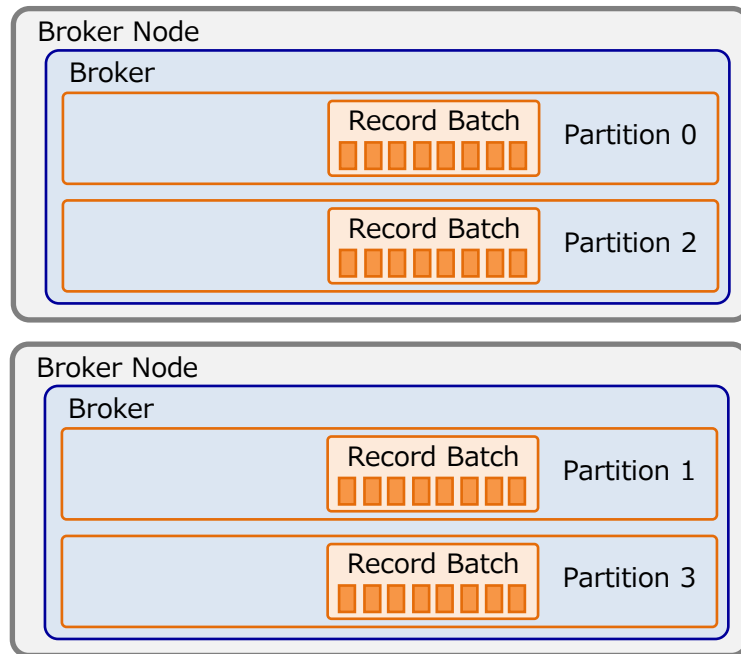
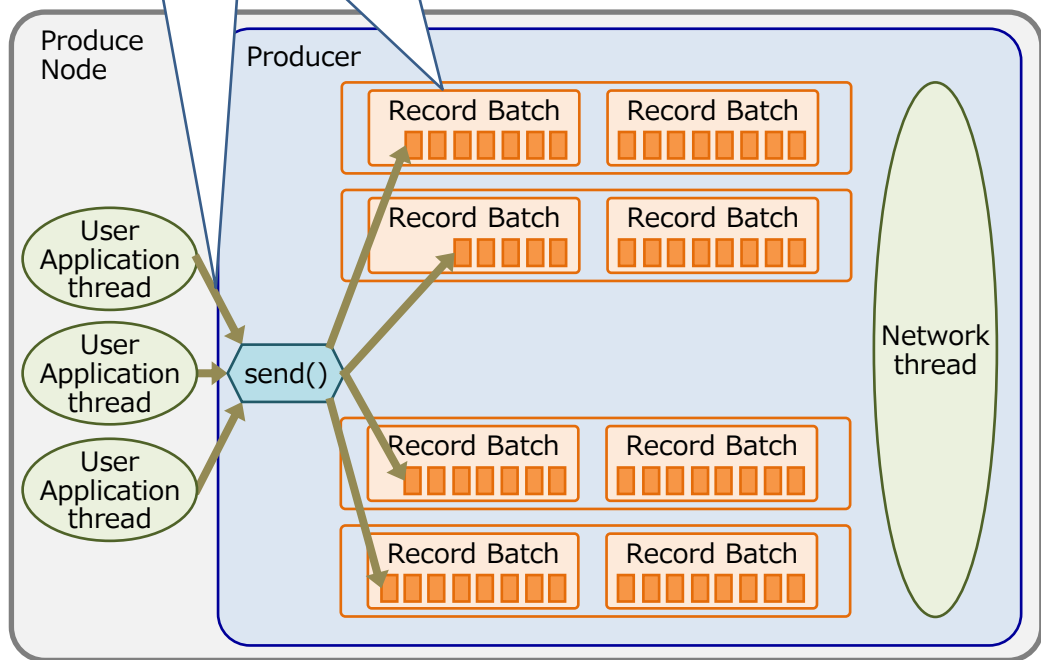
1. Recordを追加



Producerの処理の流れ (2/5)

2. Record Batchに振り分けてバッファリング
圧縮設定をした場合はRecord Batch単位で圧縮

1. Recordを追加

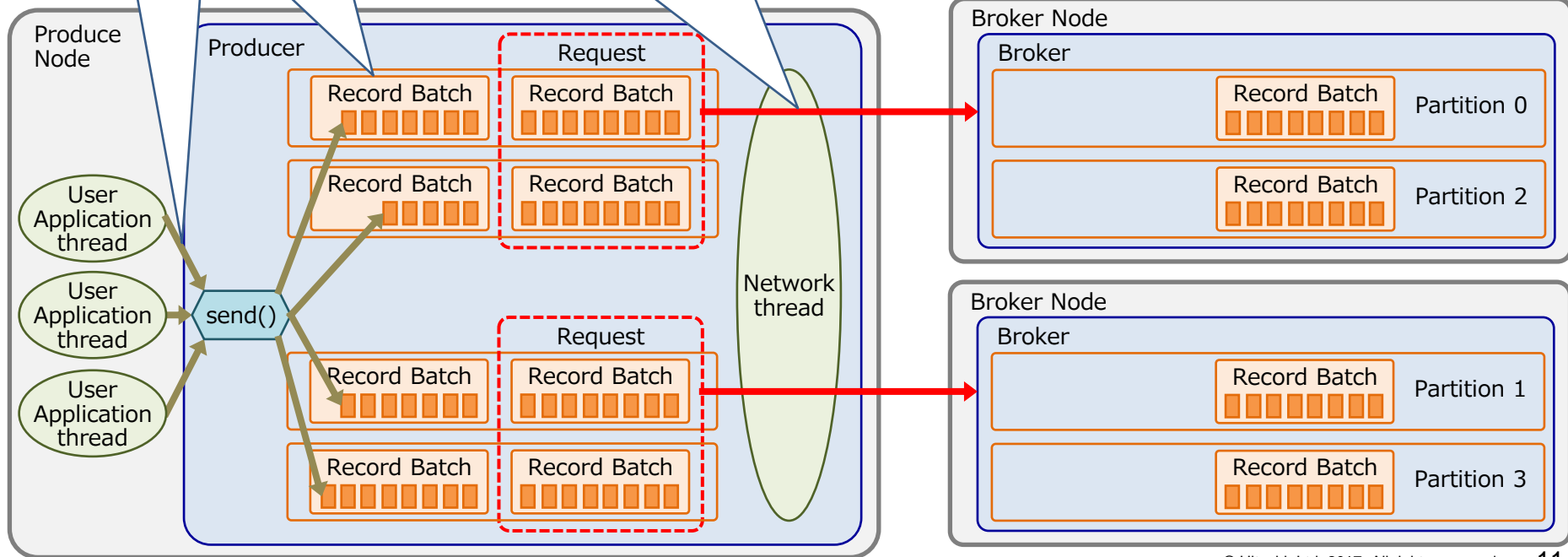


Producerの処理の流れ (3/5)

2. Record Batchに振り分けてバッファリング
圧縮設定をした場合はRecord Batch単位で圧縮

3. 複数のRecord Batchを
Broker単位でまとめて送信
(これをリクエストと呼ぶ)

1. Recordを追加



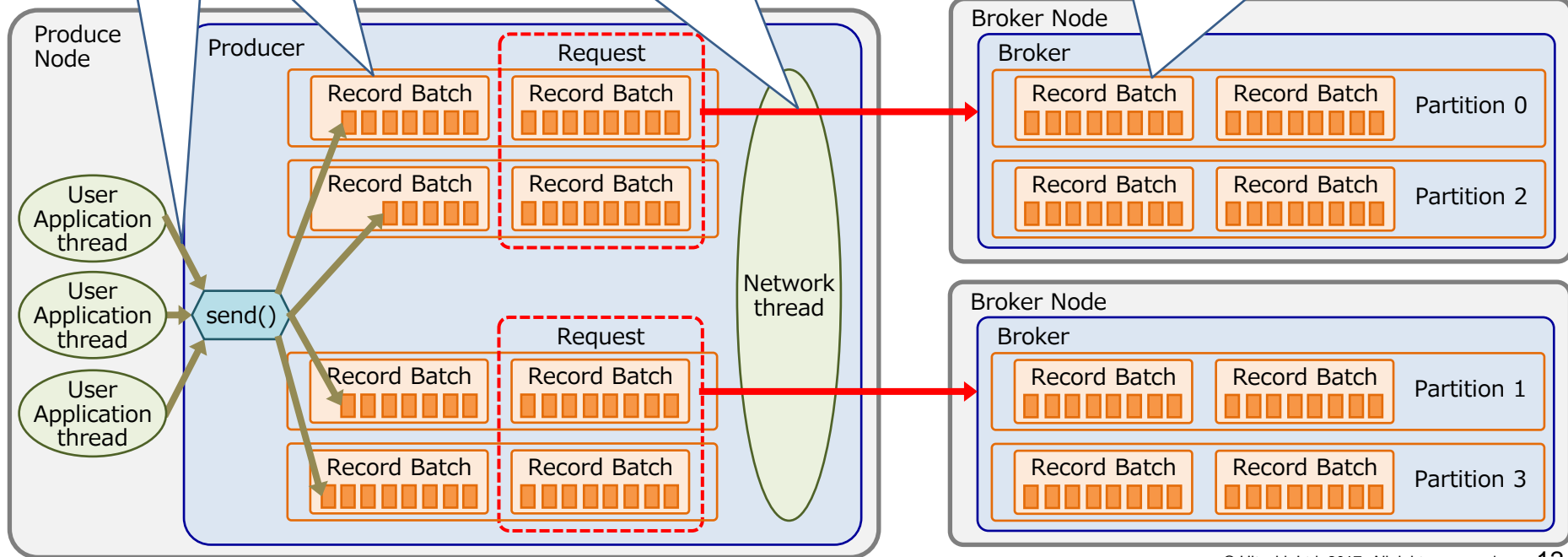
Producerの処理の流れ (4/5)

2. Record Batchに振り分けてバッファリング
圧縮設定をした場合はRecord Batch単位で圧縮

3. 複数のRecord Batchを
Broker単位でまとめて送信
(これをリクエストと呼ぶ)

4. Record Batchを対応Partitionに格納
Record Batchは未解凍のまま

1. Recordを追加



Producerの処理の流れ (5/5)

2. Record Batchに振り分けてバッファリング
圧縮設定をした場合はRecord Batch単位で圧縮

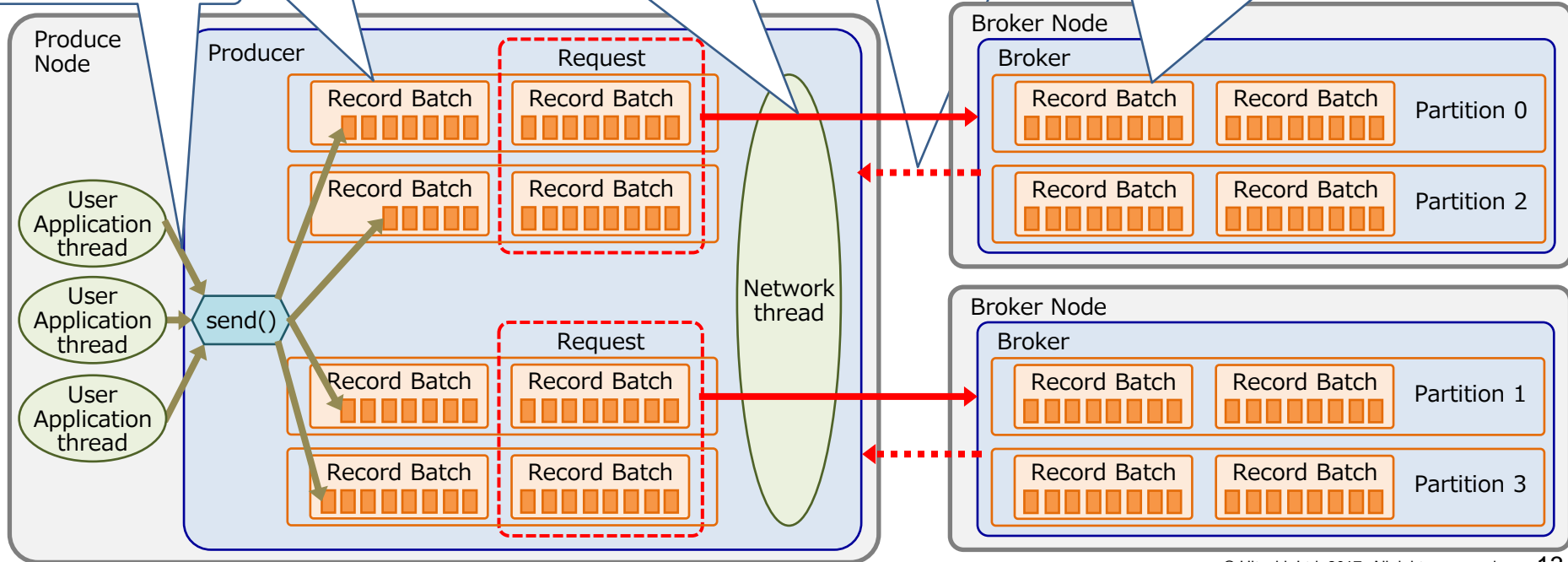
5. 指定タイミングでリクエスト完了通知 (ack) を返信

- acks=0 : ack返信なし
- acks=1 : Leader Replica書き込み完了時
- acks=all : 最小ISR(In-sync replicas)数まで複製完了時

1. Recordを追加

3. 複数のRecord Batchを
Broker単位でまとめて送信
(これをリクエストと呼ぶ)

4. Record Batchを対応Partitionに格納
Record Batchは未解凍のまま



Producerのチューニングポイントとなるパラメータとデフォルト値

Record Batch

- batch.size=16KB
- compression.type=none

メモリ使用量

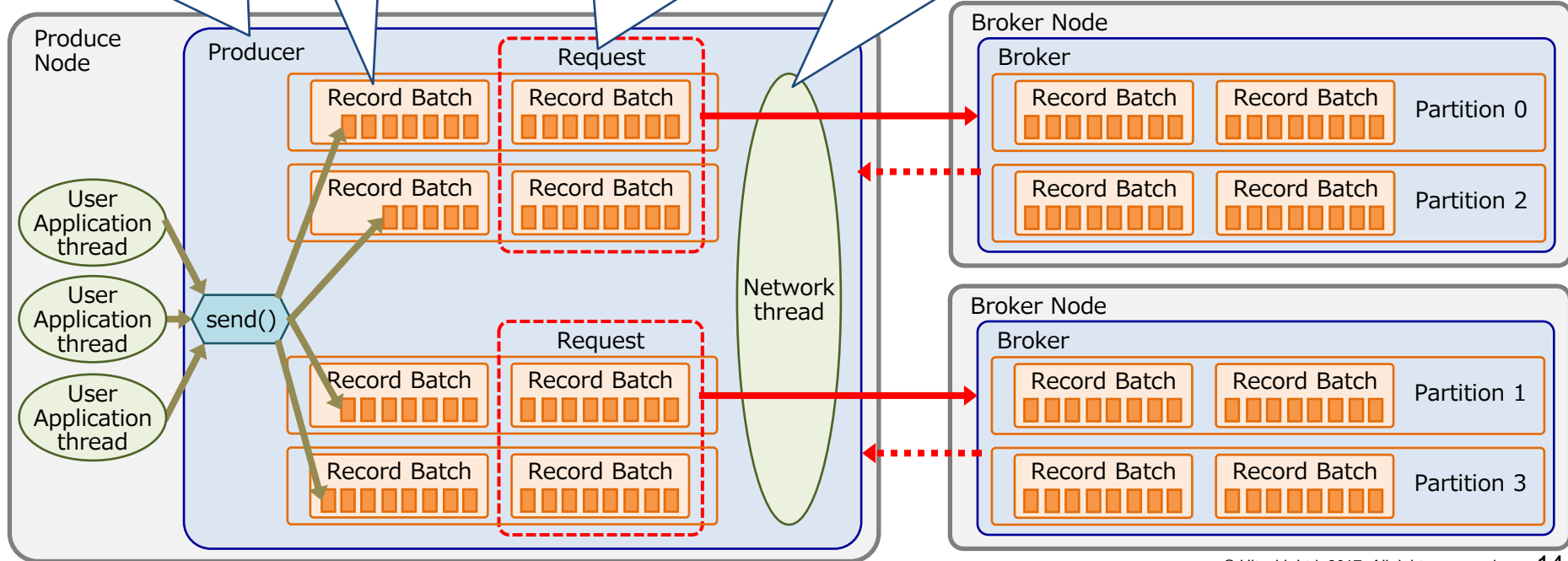
- buffer.memory=32MB

リクエスト

- max.request.size=1MB
- acks=1

リクエスト送信スレッド

- linger.ms=0ms (データ蓄積の最大待機時間)
- max.in.flight.requests.per.connection=5
- send.buffer.bytes=128KB



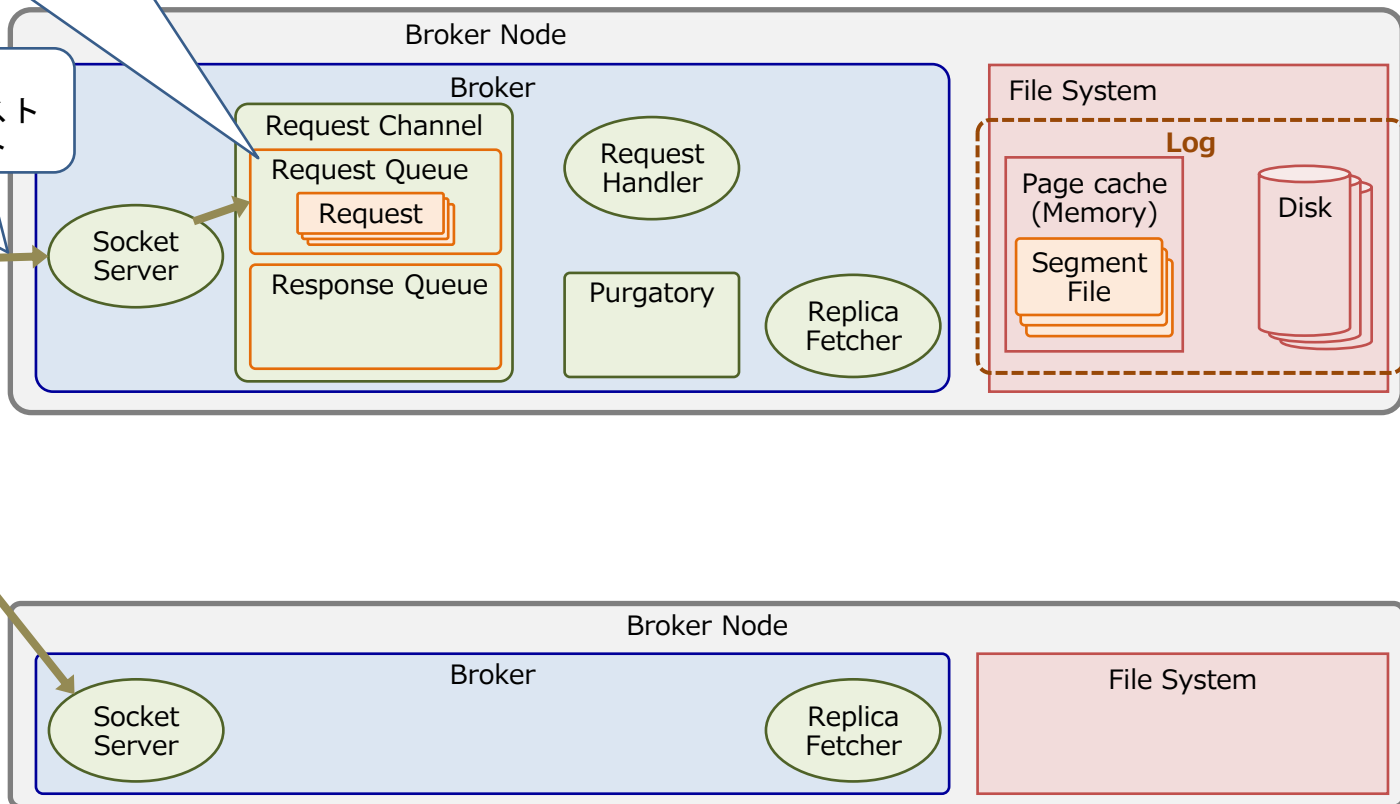
Brokerの処理の流れ [Produce (格納) / Fetch (取得) リクエスト時] (1/4)

2. リクエストをキューイング

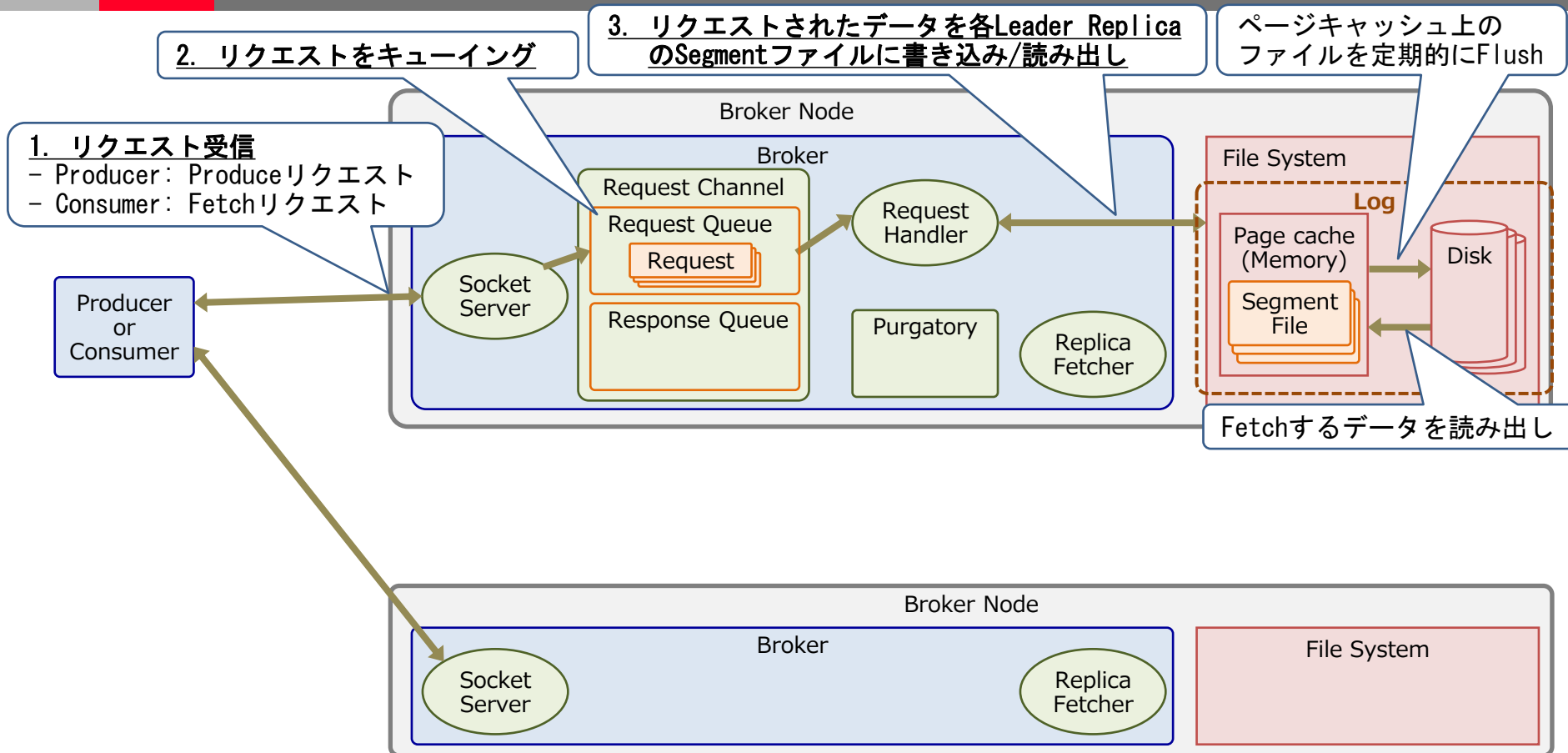
1. リクエスト受信

- Producer: Produceリクエスト
- Consumer: Fetchリクエスト

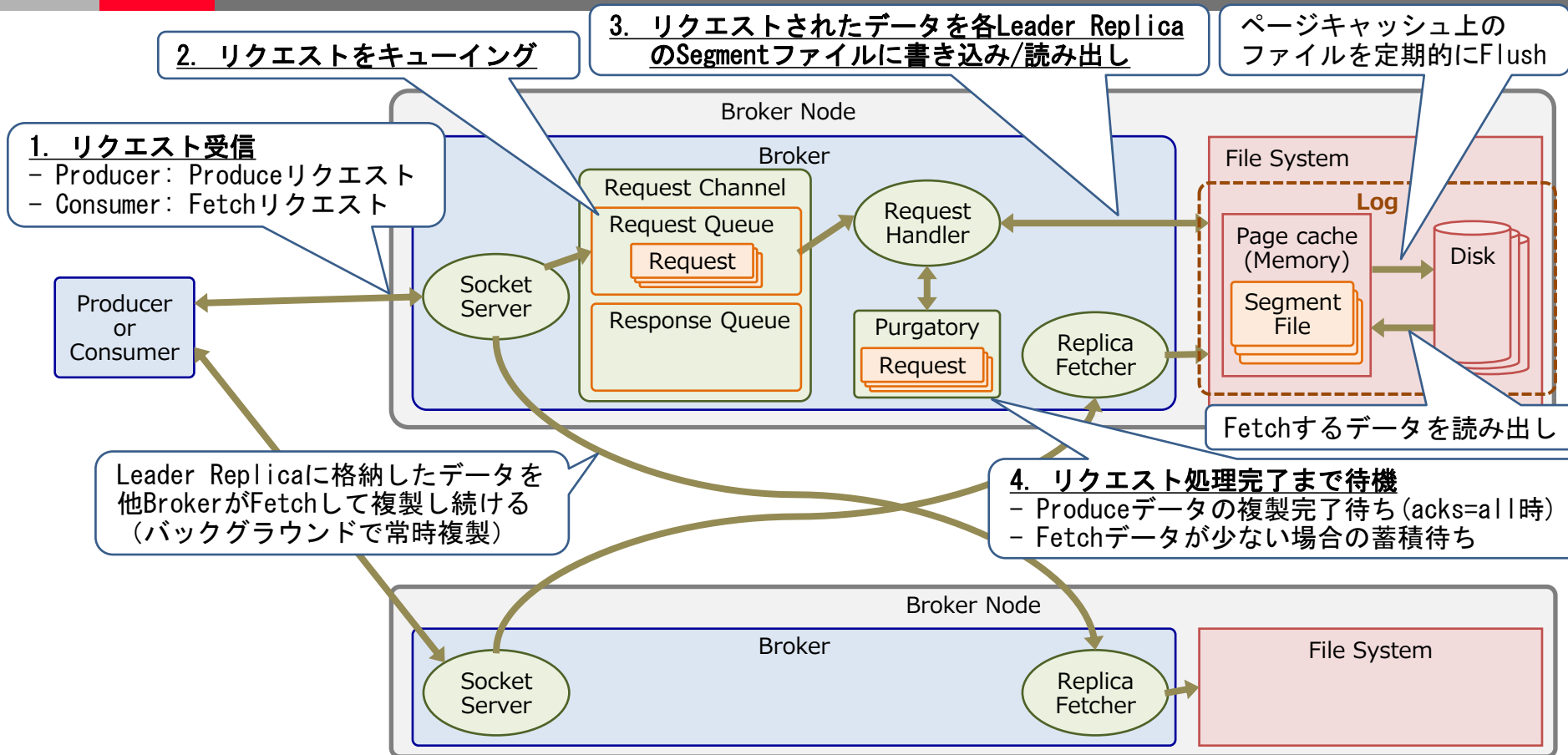
Producer
or
Consumer



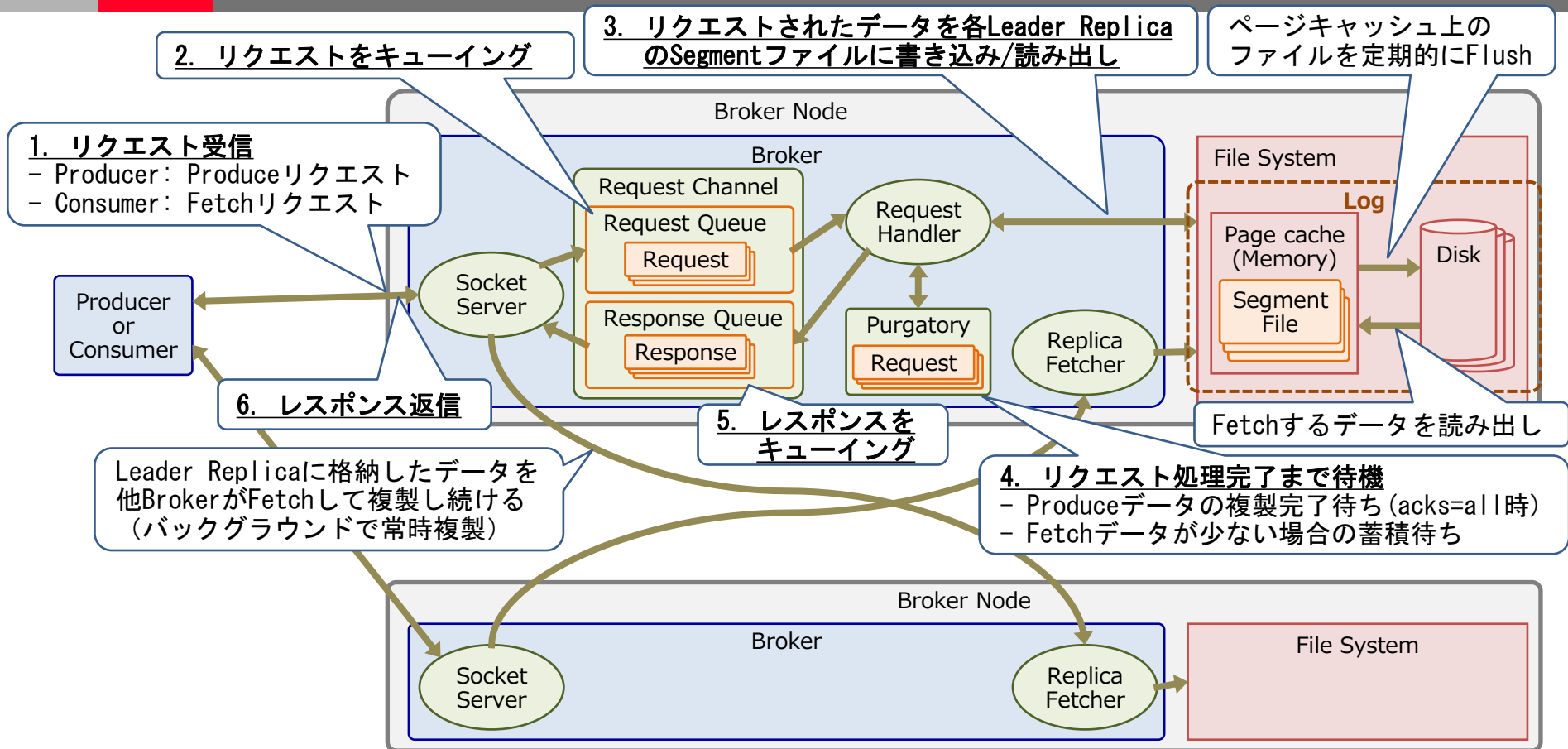
Brokerの処理の流れ [Produce (格納) / Fetch (取得) リクエスト時] (2/4)



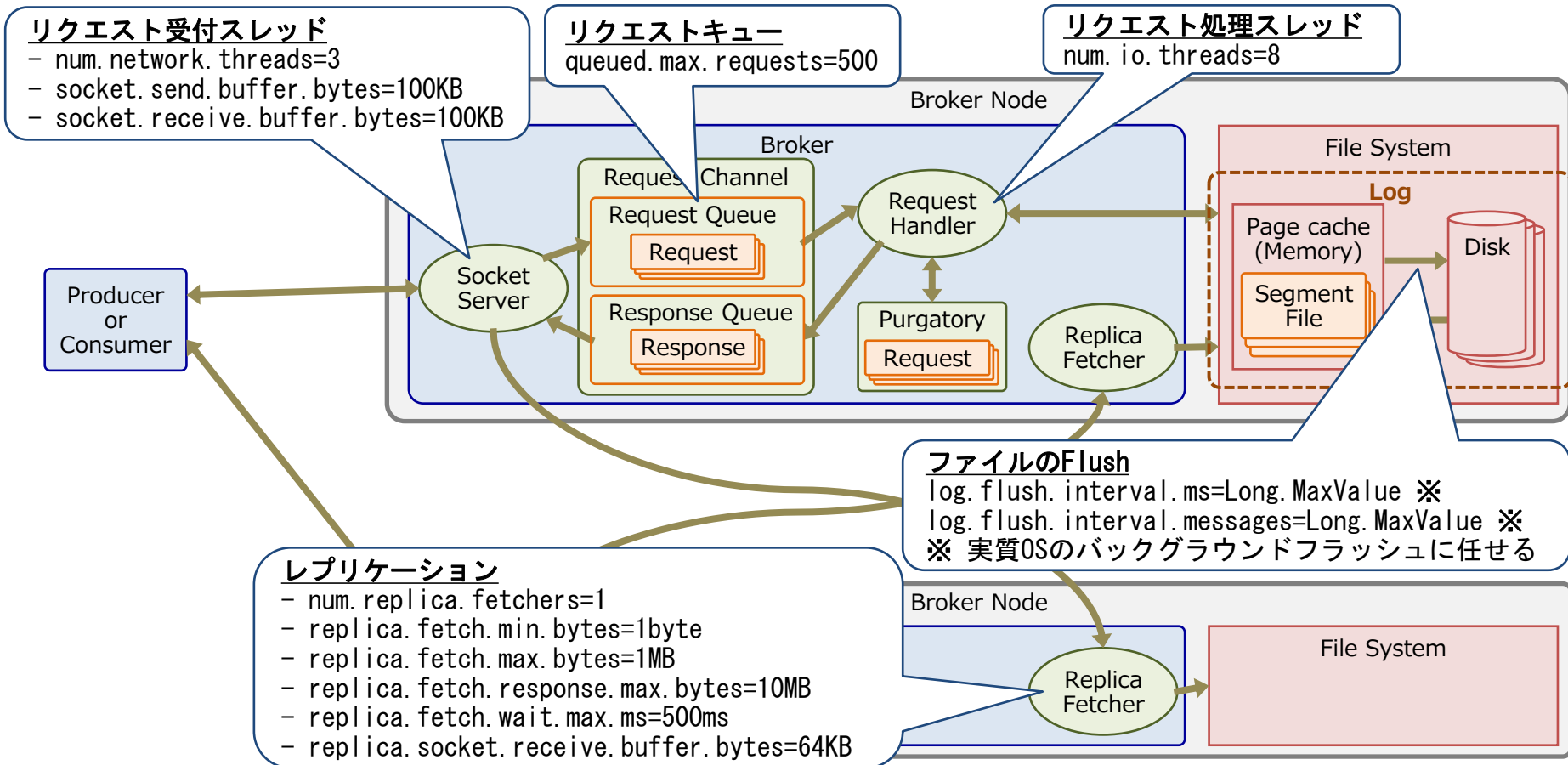
Brokerの処理の流れ [Produce (格納) / Fetch (取得) リクエスト時] (3/4)



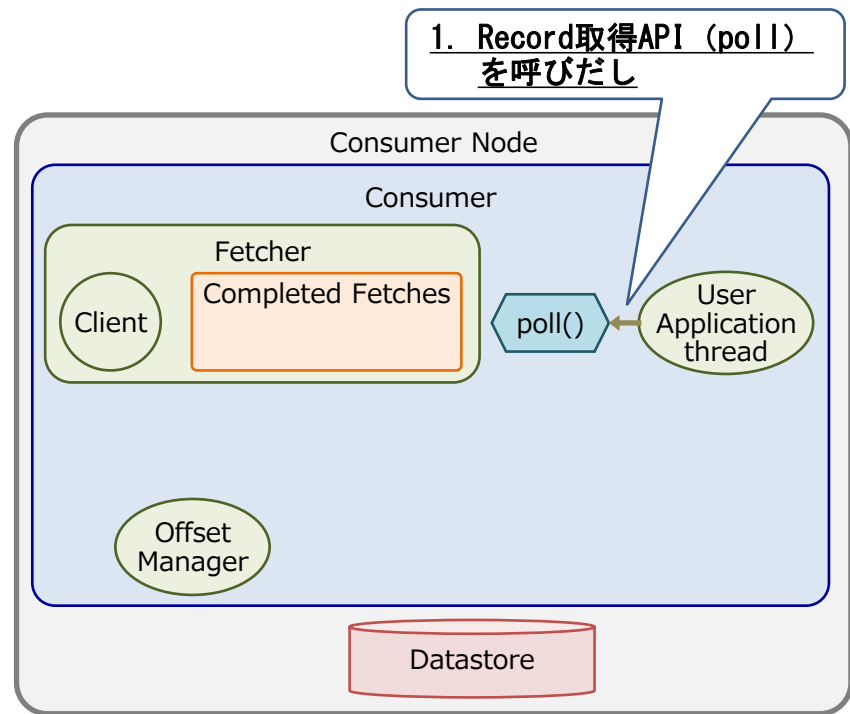
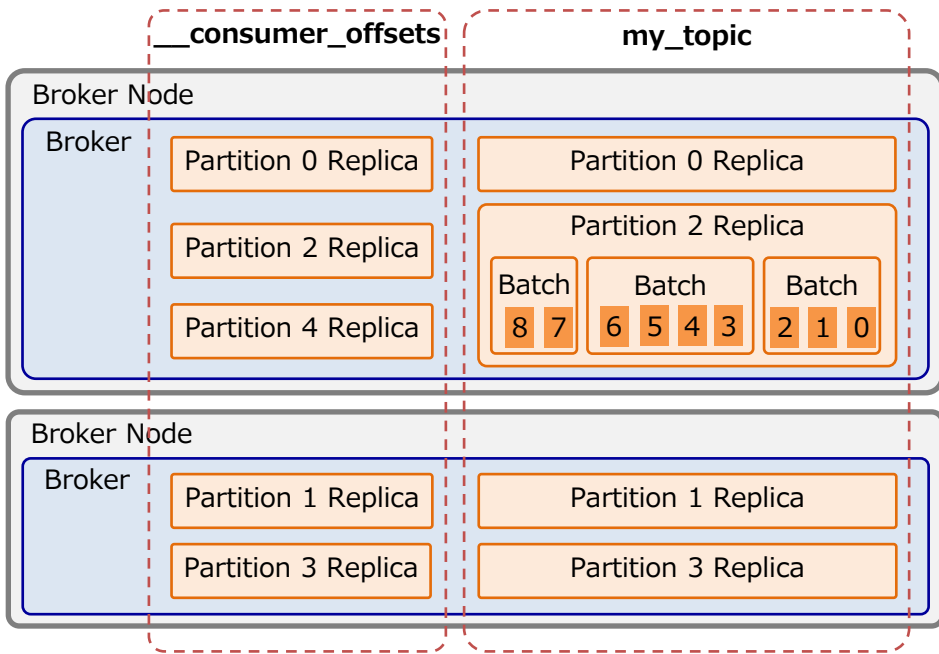
Brokerの処理の流れ [Produce (格納) / Fetch (取得) リクエスト時] (4/4)



Brokerのチューニングポイントとなるパラメータとデフォルト値



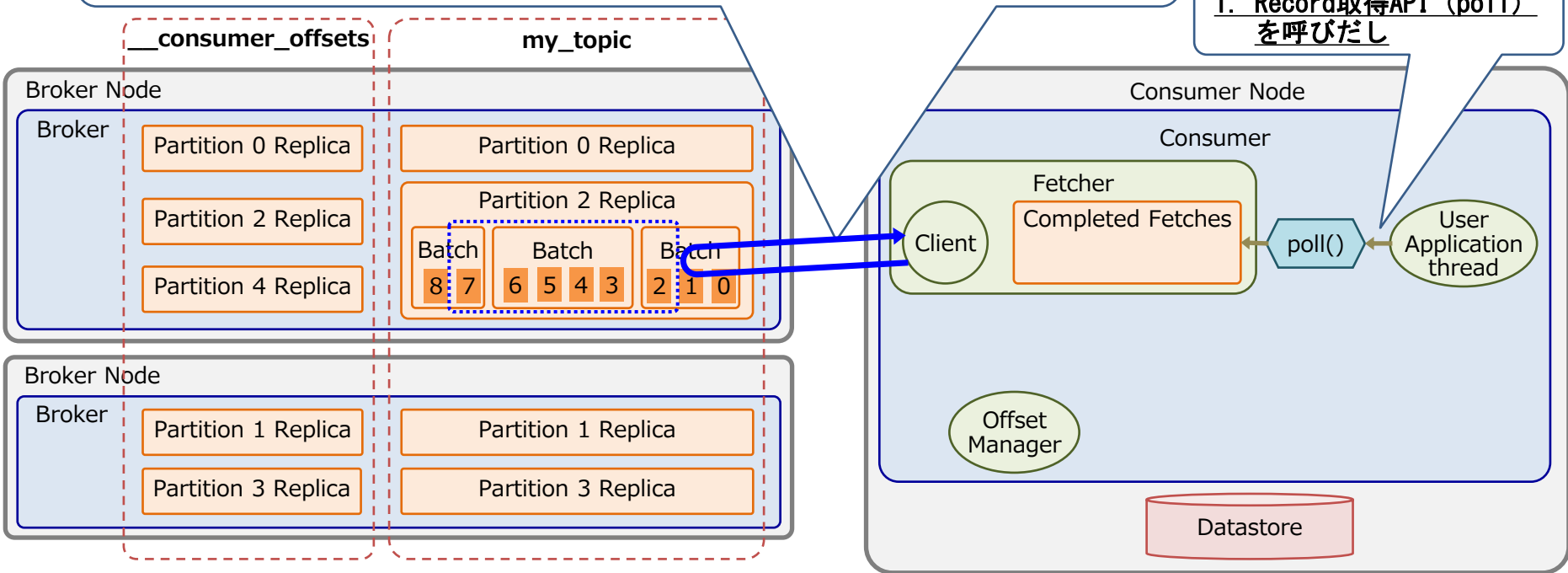
Consumerの処理の流れ(1/4)



Consumerの処理の流れ(2/4)

2. 取得対象RecordがFetcherのキューに無い場合、FetchリクエストでBrokerから取得
例 : Topic = my_topic, Partition = 2, Offset = 2-7
実際の取得単位はRecord Batch単位となるため、Offset = 0-8 のRecordを取得

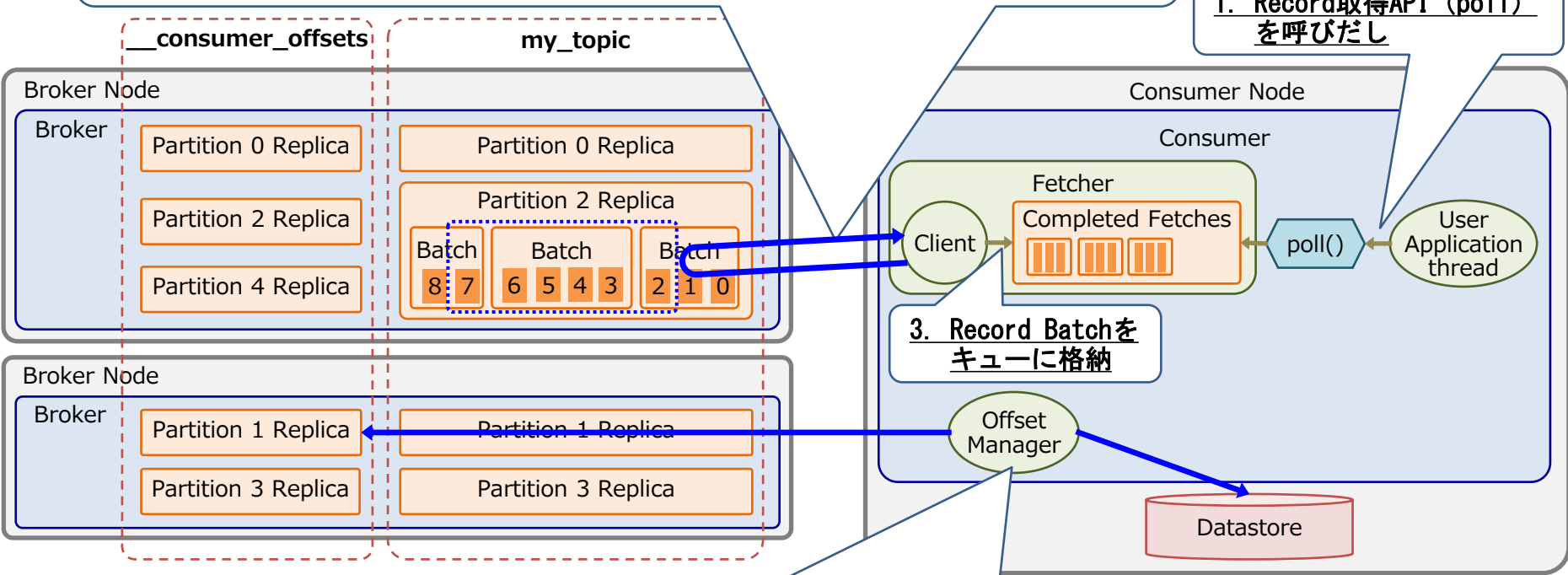
1. Record取得API (poll())
を呼びだし



Consumerの処理の流れ(3/4)

2. 取得対象RecordがFetcherのキューに無い場合、FetchリクエストでBrokerから取得
例 : Topic = my_topic, Partition = 2, Offset = 2-7
実際の取得単位はRecord Batch単位となるため、Offset = 0-8 のRecordを取得

1. Record取得API (poll())
を呼びだし

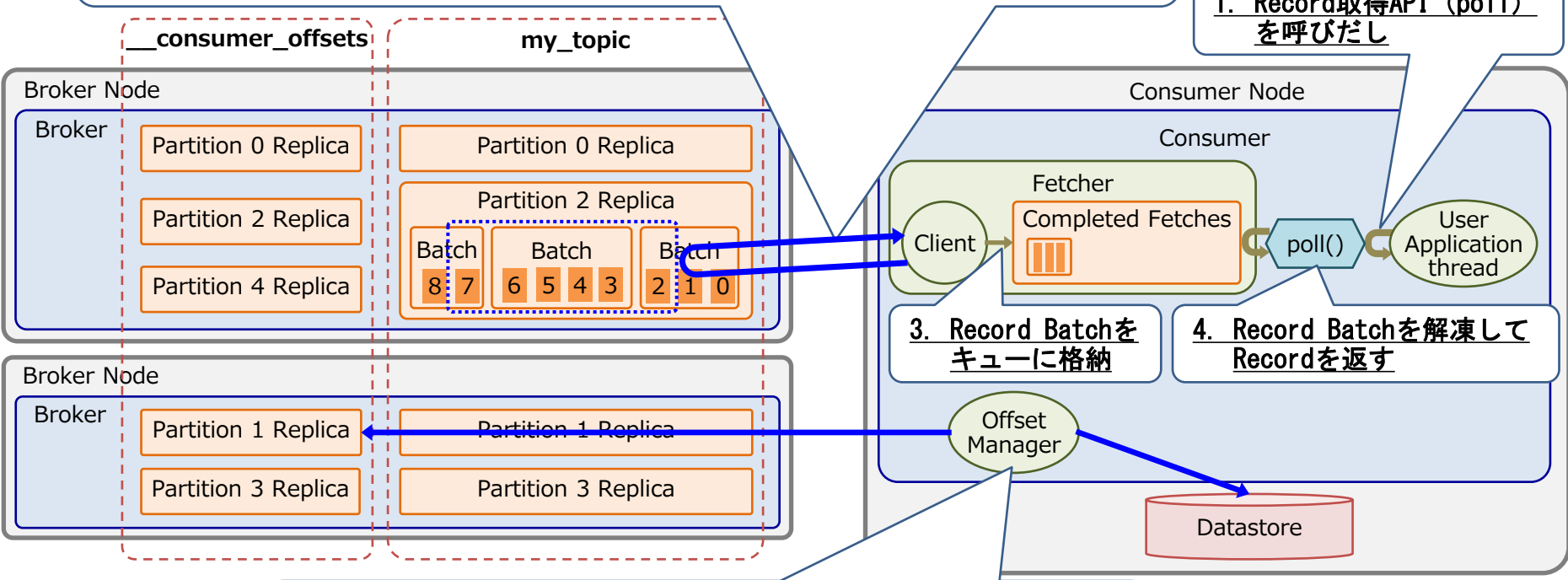


どこまで読み出したかを示すOffsetはConsumer側で管理する。
Offsetは、KafkaのTopicや任意のデータストアに保存して永続化。

Consumerの処理の流れ(4/4)

2. 取得対象RecordがFetcherのキューに無い場合、FetchリクエストでBrokerから取得
例 : Topic = my_topic, Partition = 2, Offset = 2-7
実際の取得単位はRecord Batch単位となるため、Offset = 0-8 のRecordを取得

1. Record取得API (poll())
を呼びだし

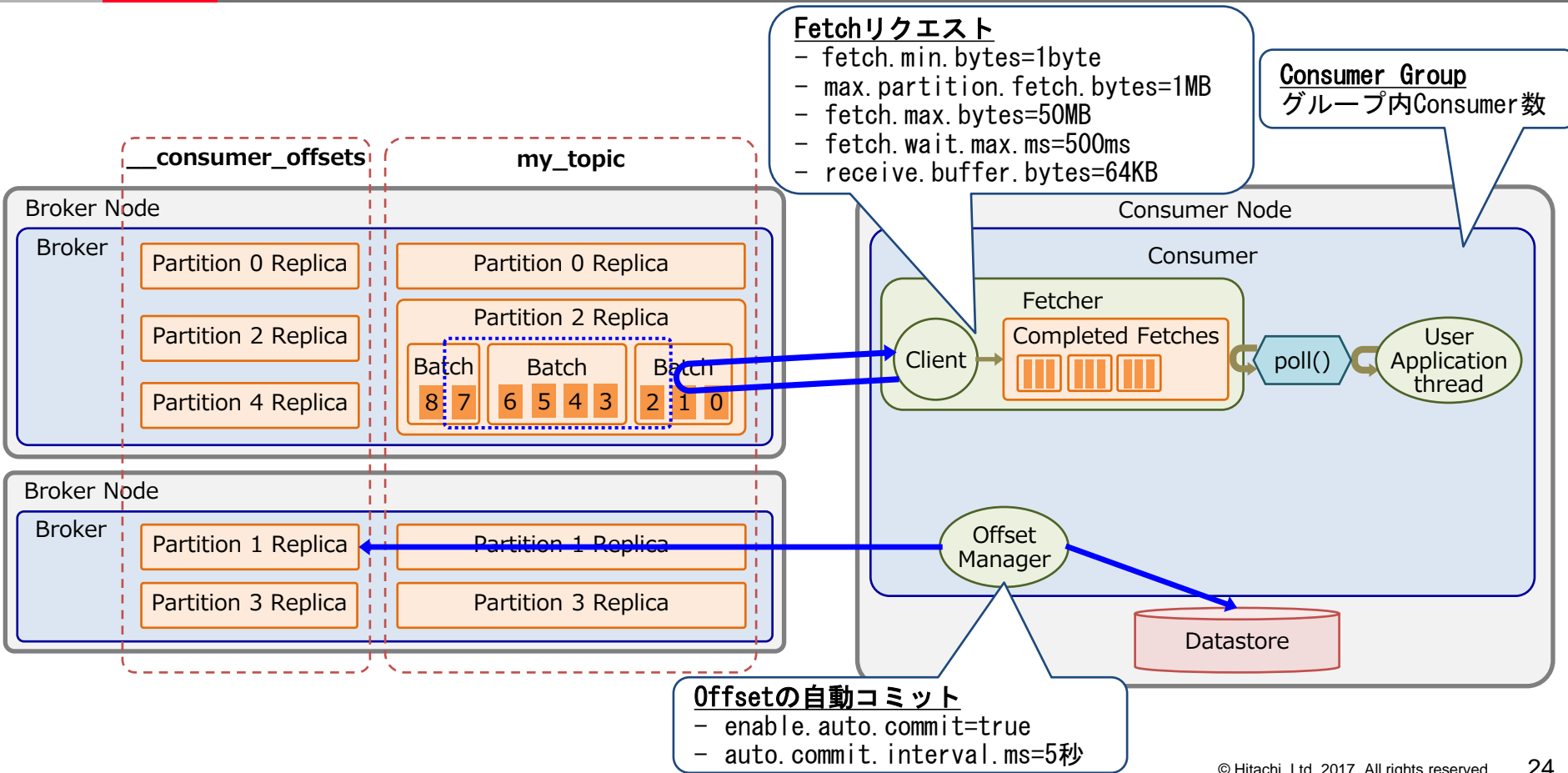


3. Record Batchを
キューに格納

4. Record Batchを解凍して
Recordを返す

どこまで読み出したかを示すOffsetはConsumer側で管理する。
Offsetは、KafkaのTopicや任意のデータストアに保存して永続化。

Consumerのチューニングポイントとなるパラメータとデフォルト値



3. 性能検証の概要

• 背景

- Kafkaは非同期レプリケーションであれば、ディスク/ネットワーク性能の上限に近いスループットを容易に出せる
- しかし、同期レプリケーション時の性能については情報が少ない
 - 本番環境ではデータを保護するため同期レプリケーションが求められる

• 検証内容

- 同期レプリケーションでどこまでスループットを出せるか確認する
 - Producerノードで1KBのRecordを生成・送信し、1Consumer Groupが受信し続け、600秒間の平均スループットを測定
 - スループットが増加するとレイテンシも増加するが、チューニングでは考慮しない
- 同期レプリケーションの設定
 - Replication factor = 3
 - 最小 In Sync Replica 数 = 2
 - acks = all
 - Producerはデータが最小ISR数(2個)以上のReplicaに同期されたことを確認する(Leader Replica + 1個以上のFollower Replicaに複製されたことを確認)

各ノードの最大スループット

- ネットワーク: 1,170MB/s (送信/受信共に)
- ディスク: 1,200MB/s (書き込み/読み出しの合計)

仮想マシン

- OS: CentOS 6.7
- CPU: 2 core
- Memory : 16 GB
- Disk: 160GB * 1



ZooKeeper

10 Gbps SW: Brocade VDX6740

1 Gbps LAN

10 Gbps LAN (約 **1,170 MB/s**)
全二重通信が可能

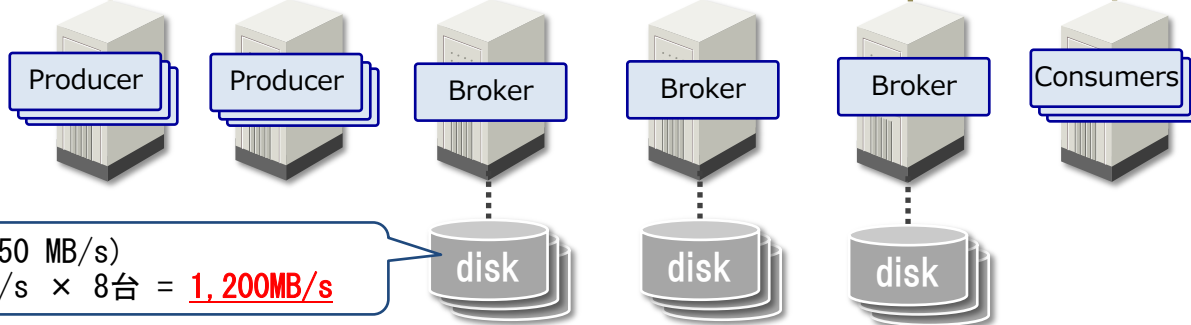


ソフトウェア

- Apache Kafka 0.11.0
- ZooKeeper 3.4.10
- Java 1.8.0.121

物理マシン: HA8000/TS20AN × 6台

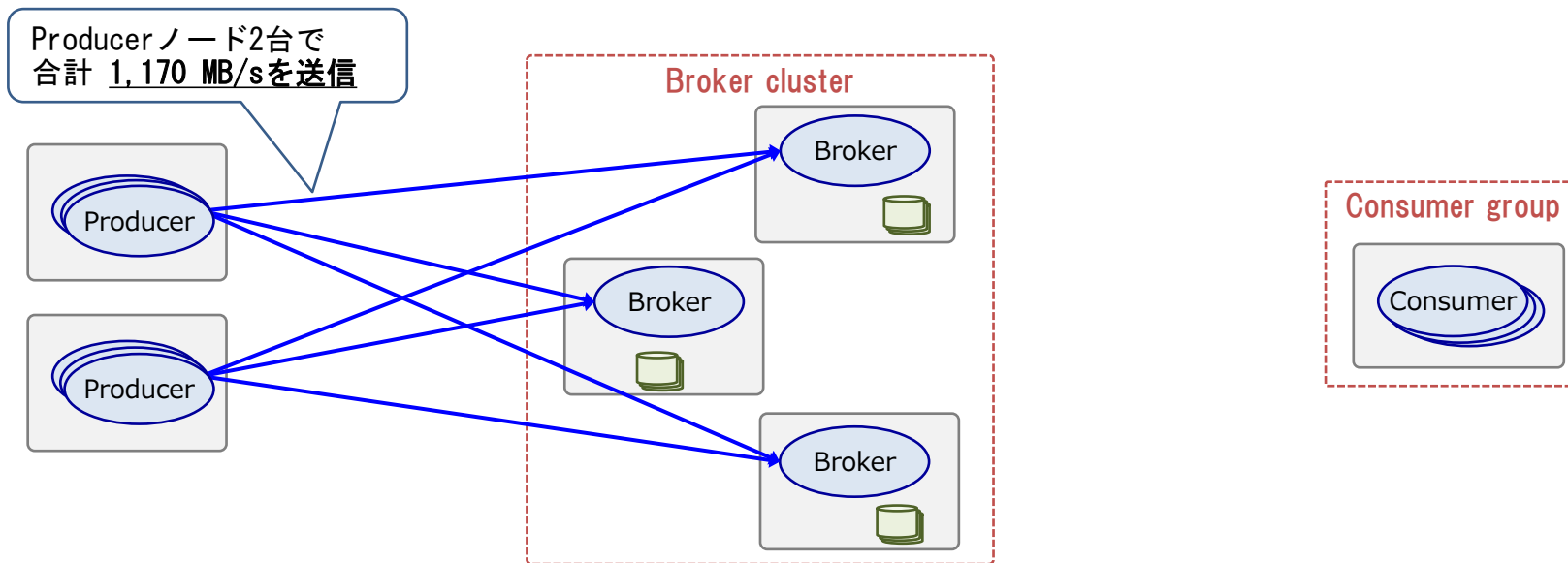
- OS: RHEL 6.7
- CPU: 40 core
- Memory : 384 GB
- Disk: 1,200GB (SAS 10,000 rpm) * 10
(OS用: 2台でRAID1, Broker用: JBODで8台)



SAS 10,000 rpm (約 150 MB/s)
1ノードあたり 150 MB/s × 8台 = **1,200MB/s**

理論上の最大スループット (1/3)

- Producerの送信/Consumerの受信スループットが釣り合う最大スループットを計算



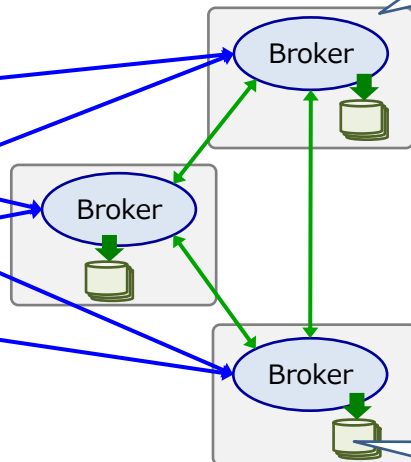
- Producerの送信/Consumerの受信スループットが釣り合う最大スループットを計算

- Producerノード2台から合計 390 MB/sで受信
- 受信データを別のBroker2台へそれぞれ 390 MB/sで送信
- 別のBroker2台からデータをそれぞれ 390 MB/sで受信

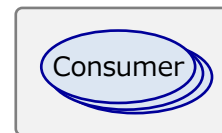
Producerノード2台で
合計 1,170 MB/sを送信



Broker cluster



Consumer group

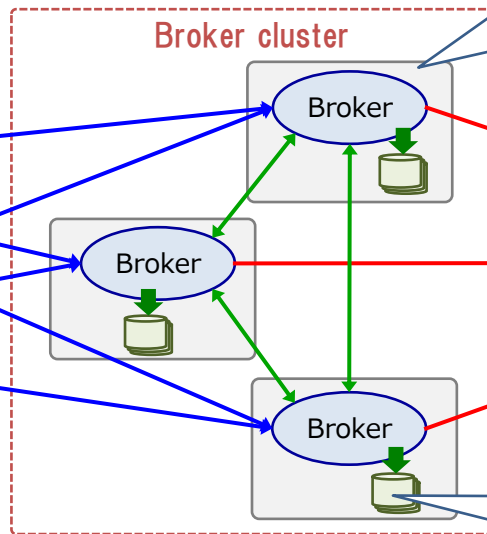
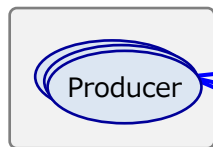


- 受信データ 1,170 MB/sをディスク書き込み

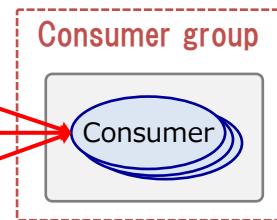
- Producerの送信/Consumerの受信スループットが釣り合う最大スループットを計算

- Producerノード2台から合計 390 MB/sで受信
 - 受信データを別のBroker2台へそれぞれ 390 MB/sで送信
 - 別のBroker2台からデータをそれぞれ 390 MB/sで受信
 - Consumerノードに 390 MB/sで受信
- ⇒ 合計で **受信 1,170 MB/s**, **送信 1,170 MB/s**

Producerノード2台で
合計 1,170 MB/sを送信



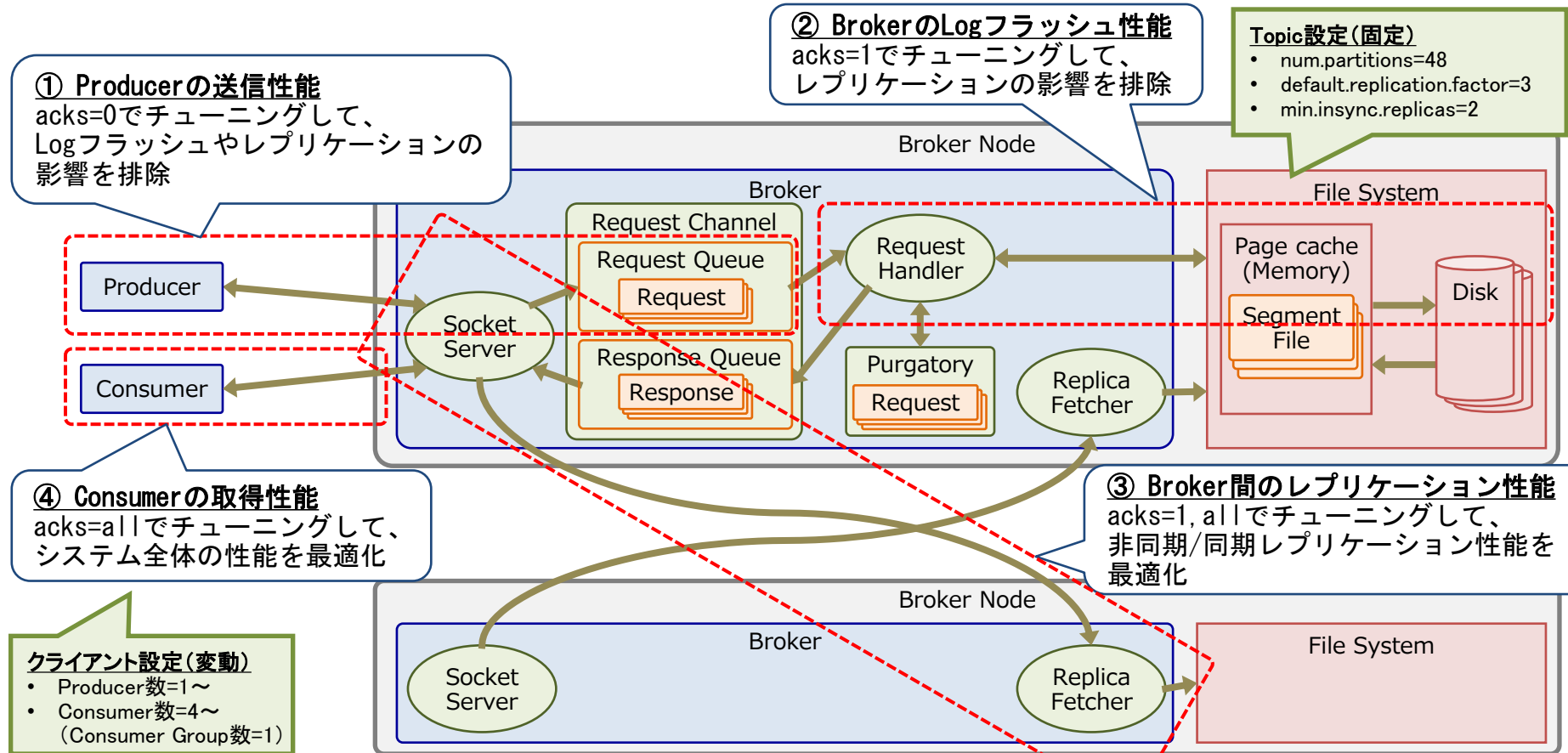
Consumerノードは
合計 1,170 MB/sを受信



- 受信データ 1,170 MB/sをディスク書き込み
- 送信データはすぐに読み出すためページキャッシュ上にある (ディスク読み出しなし)

Brokerの送受信とConsumerの受信がボトルネック: システム全体の最大スループットは 1,170 MB/s

検証範囲を4分割して順番にチューニング

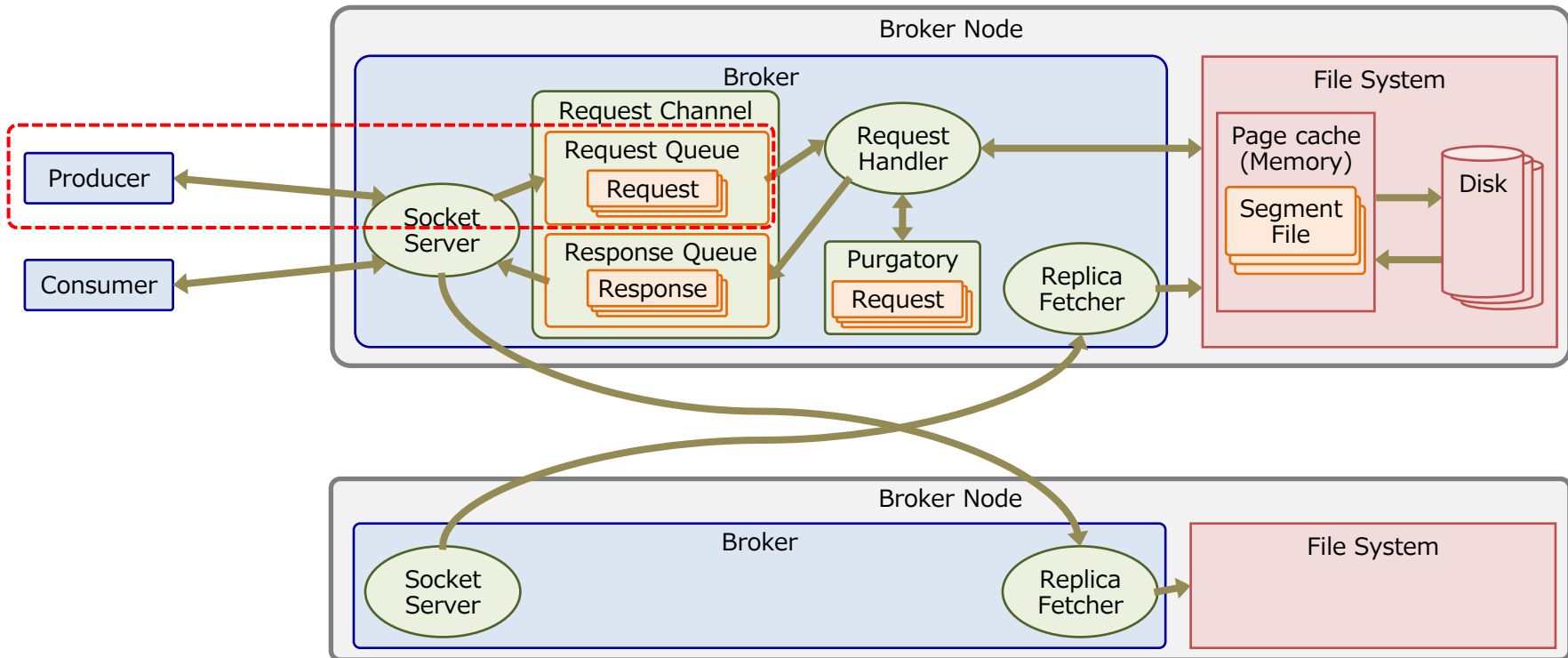


4. 検証結果と考察

① Producerの送信性能のチューニング

- 目的: acks=0 で Produceスループットを最大化する

- Brokerに十分な量のデータを入力するため、Producer送信スループット1,170MB/sをめざす



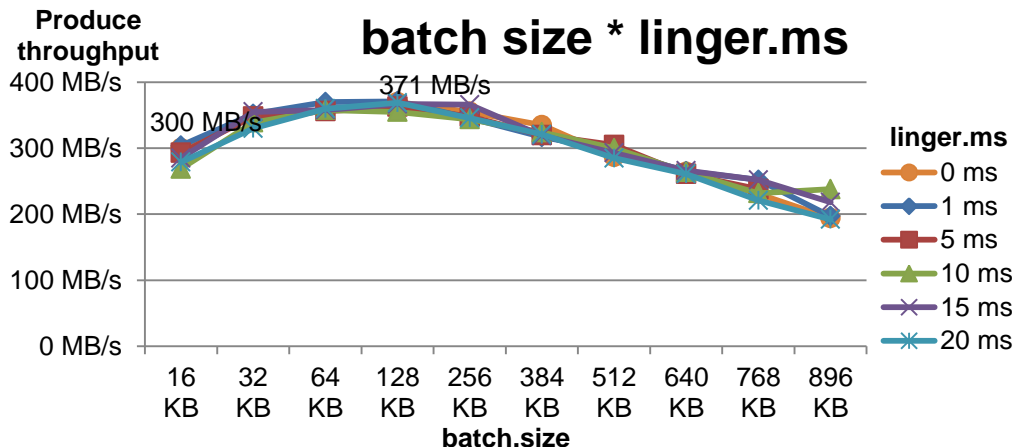
① Producerの送信性能のチューニング

- 送信するRecord Batchサイズは、最大サイズ(batch.size)と蓄積待機時間(linger.ms)で決まるため、これを同時にチューニング
 - Broker側のLogフラッシュやレプリケーションの影響を排除するため acks=0 で測定
 - Producer用ノード1台かつProducer1個で測定

Producerのパラメータ設定

- acks = 0
- buffer.memory = 32MB->1GB
- max.request.size = 1MB->最大値
- max.in.flight.requests.per.connection = 5->最大値
- max.block.ms = 60秒->最大値
- retries = 0回->最大値
- request.timeout.ms = 30秒->最大値

※ 最大値： IntegerまたはLong型の最大値を設定

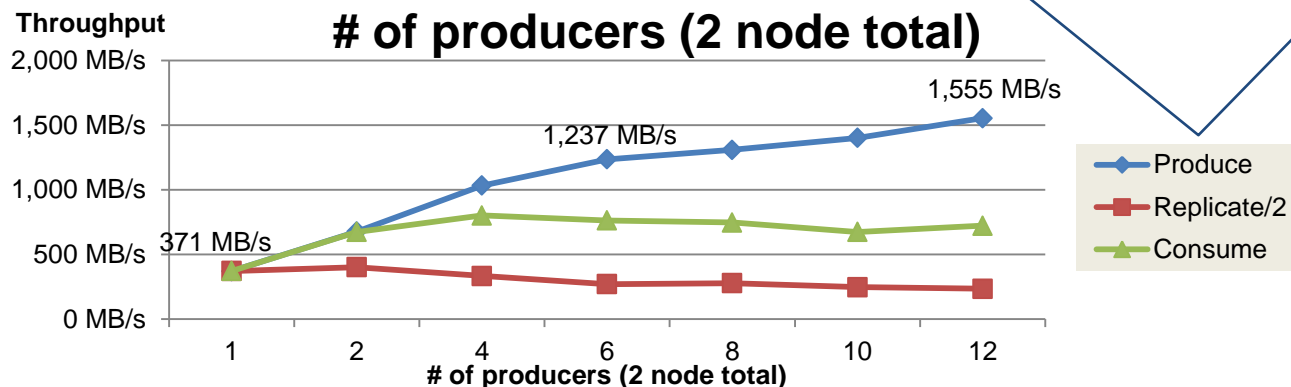


- ◆ batch.sizeのチューニングは効果大きいですが、linger.msの影響は少なかった
- ◆ batch.size=128KB, linger.ms=1msのとき、Produceスループットは 300 -> 371MB/s に向上

① Producerの送信性能のチューニング

• Producer用ノード2台でProducer数を2の倍数で増やして測定

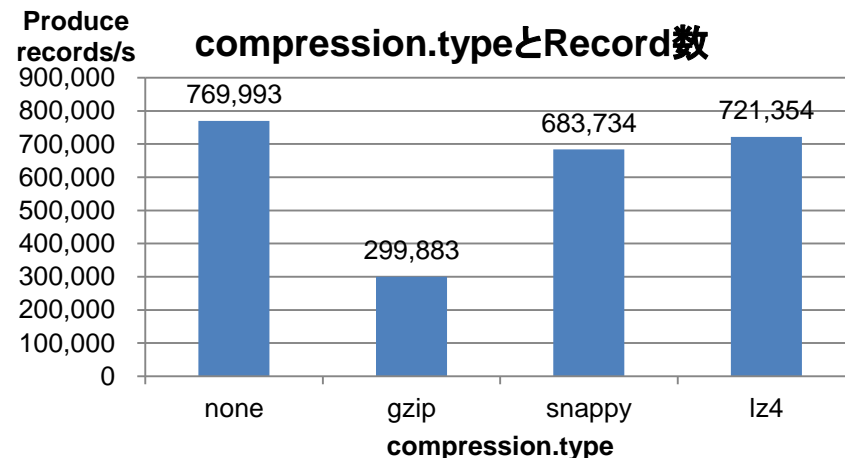
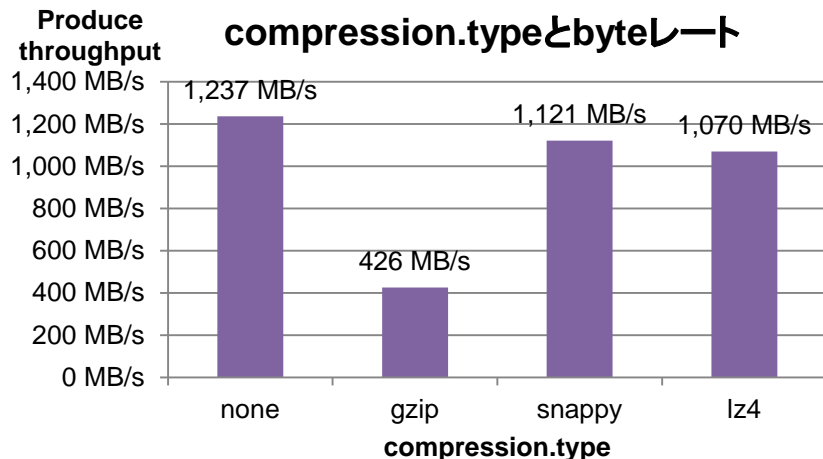
- Produce: Producerの送信スループット
- Replicate/2: Broker間レプリケーションのスループット
(同じデータを他の2台に複製するため、1/2の値を表示)
- Consume: Consumerの受信スループット



Producer数=6 でProduceスループットが理論値である 1,170 MB/sを超えた
これは、Replicateスループットが遅延した分のネットワーク帯域を使用できるためと考えられる

① Producerの送信性能のチューニング

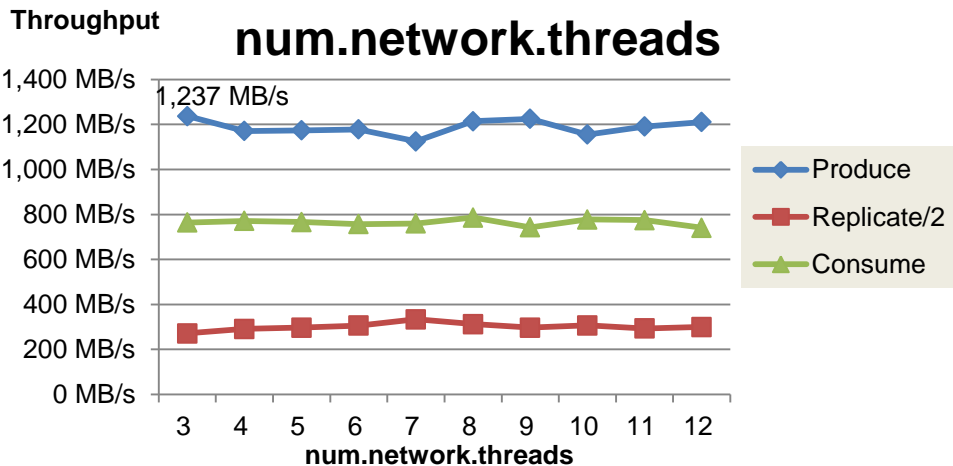
- ProducerでRecord Batchを圧縮して測定
 - 4種類の圧縮アルゴリズム(compression.type)を検証
 - 圧縮率により送信レコード数が増えるため、ByteレートとRecordレートの両方を測定



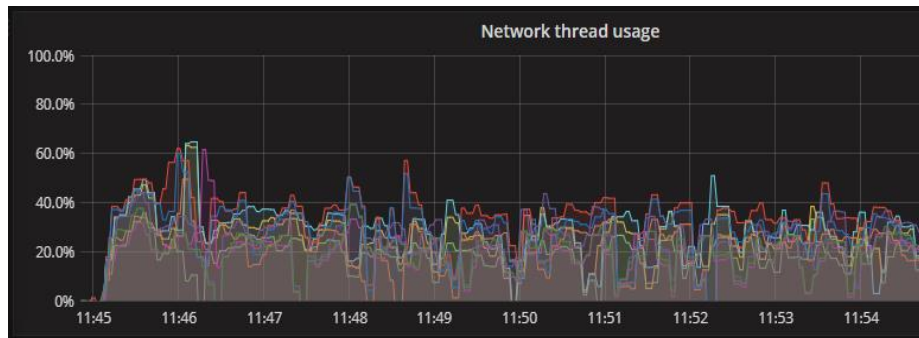
- ◆ gzipは圧縮率が高い分Byteレートは低く、snappyとlz4は圧縮率が低い分Byteレートは高い傾向
- ◆ ただしRecordレートを比べると圧縮なし(none)が最も高かった

① Producerの送信性能のチューニング

- 受信側 Broker の Socket Server のネットワークスレッド数 (num.network.threads)を増やして測定



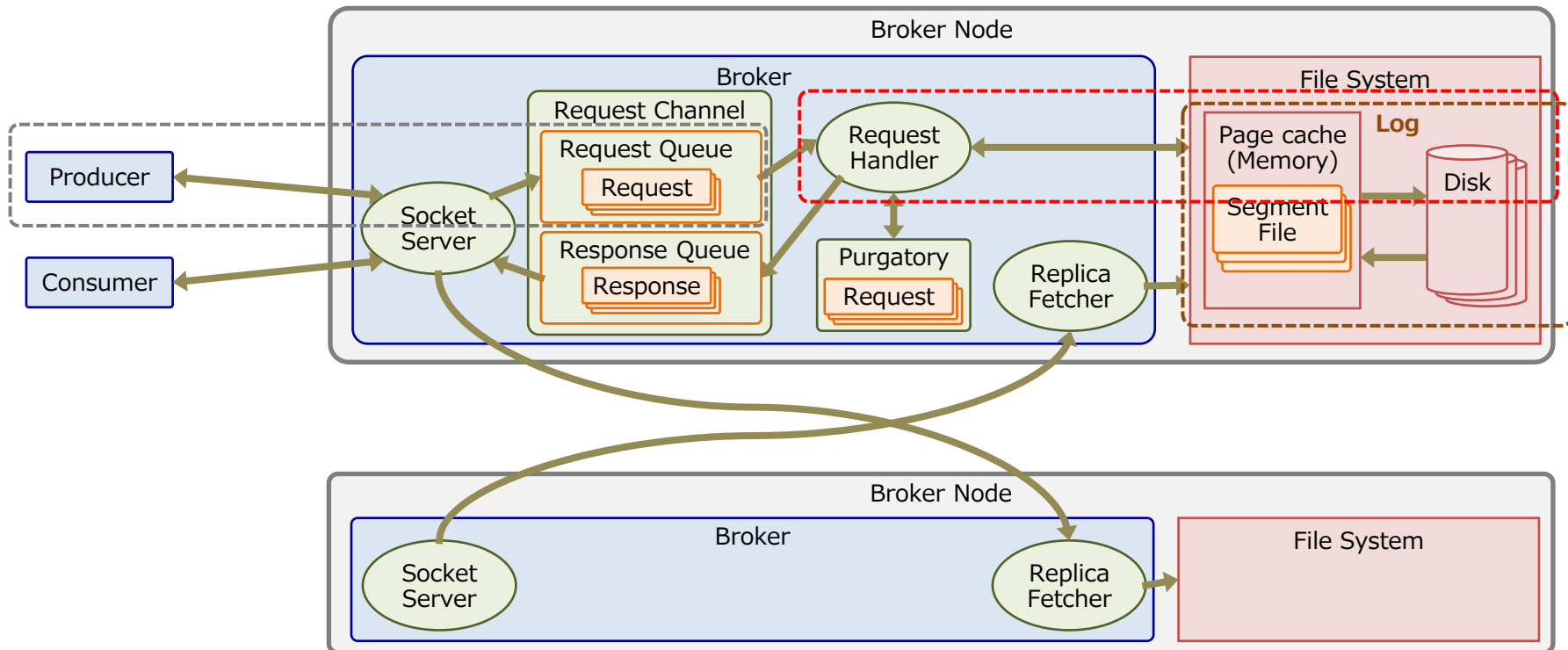
ネットワークスレッドのCPU使用率
(3 num.network.threads × 3 broker = 9 threads)



- ◆ ネットワークスレッド数を増やしてもProduceスループットはほぼ横ばいであった
- ◆ ネットワークスレッドのCPU使用率は40%程度のためボトルネックではない

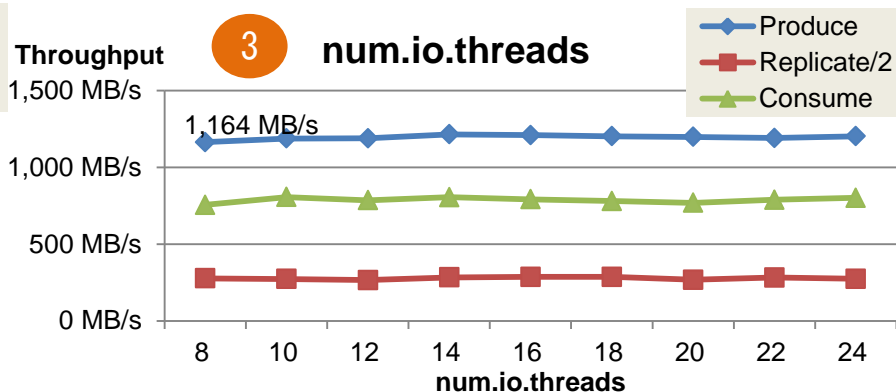
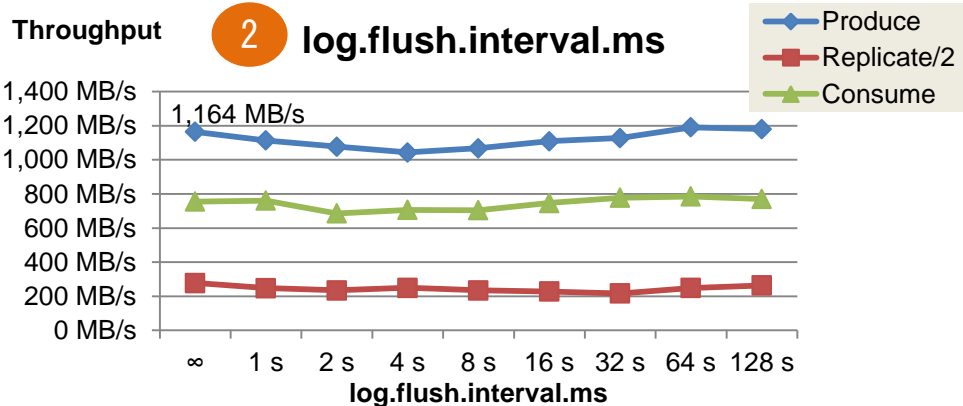
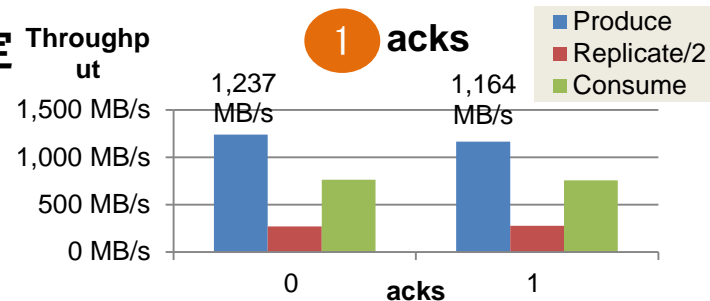
② BrokerのLogフラッシュ性能のチューニング

- 目的: acks=1 で Produceスループットを最大化する
 - Logフラッシュまで含めた性能を最適化する



② BrokerのLogフラッシュ性能のチューニング

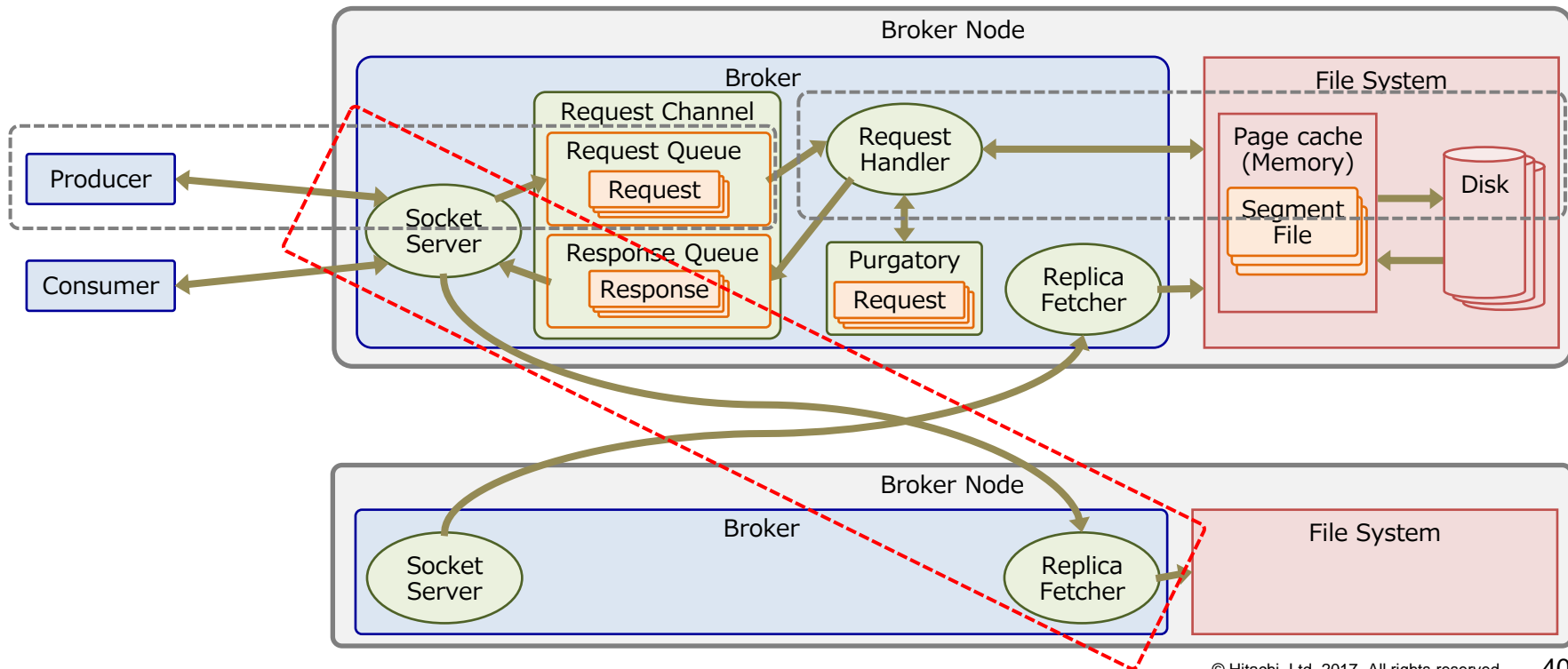
- ① Logフラッシュを含む性能を測定するためacks=1に設定
- ② Logフラッシュ間隔をチューニング
 - デフォルト設定ではOSのフラッシュに任せている
- ③ リクエスト処理スレッド数をチューニング



- ◆ acks=0→1でディスクI/Oがボトルネックとなり、Produceスループットは1,237 → 1,164MB/sに低下
- ◆ log.flush.interval.ms / num.io.threads のチューニングはほぼ効果がなかった

③ Broker間のレプリケーション性能のチューニング

- 目的: acks=all で Replicateスループットを最大化する
 - Produce/Replicateスループットの変化を見るため、まずはacks=1でチューニングを行い、最後にacks=allに変更する

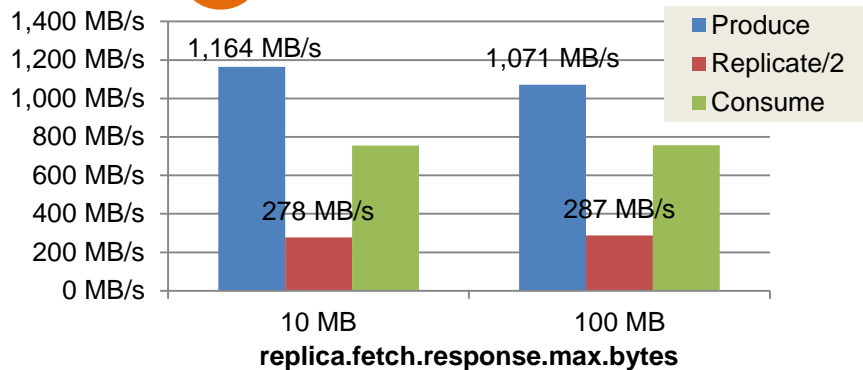


③ Broker間のレプリケーション性能のチューニング

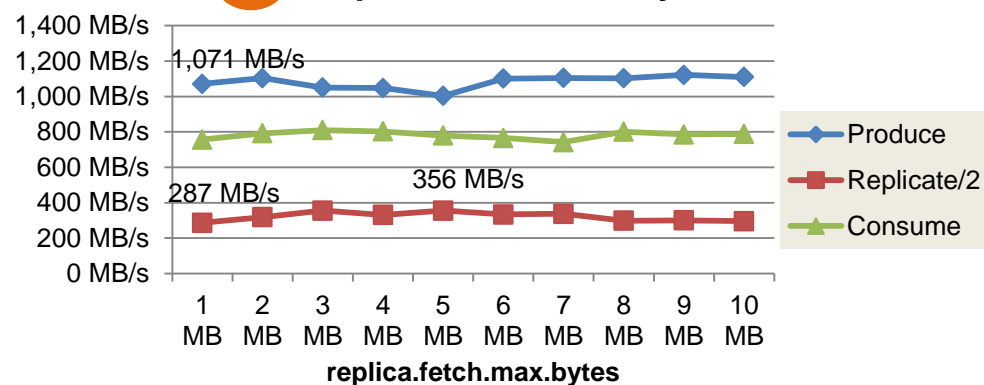
- Replicateスループットを最大化するため、Replica fetcherが送信するFetchリクエストの設定をチューニング

- Fetchリクエストで取得する最大データサイズを変更
- Fetchリクエストで取得するPartition単位の最大データサイズを増やす

Throughput ① replica.fetch.response.max.bytes



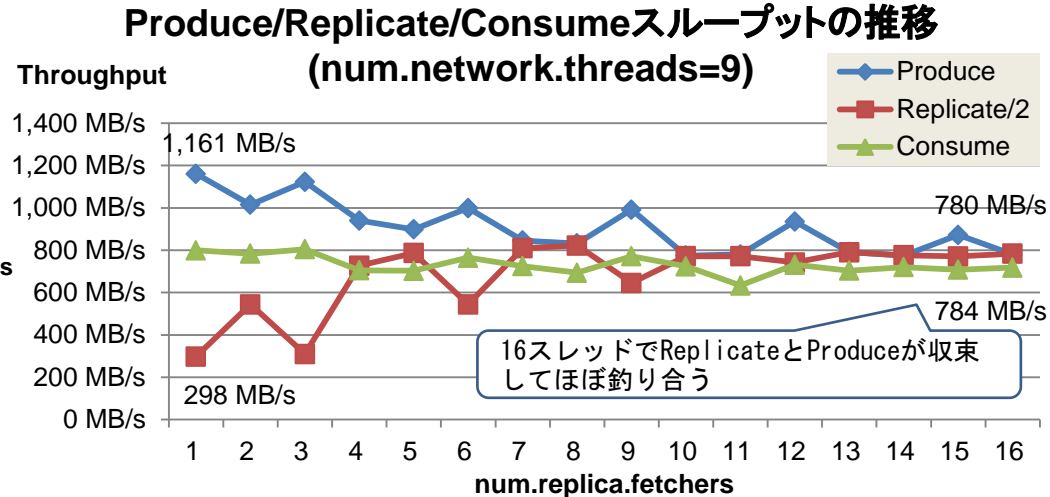
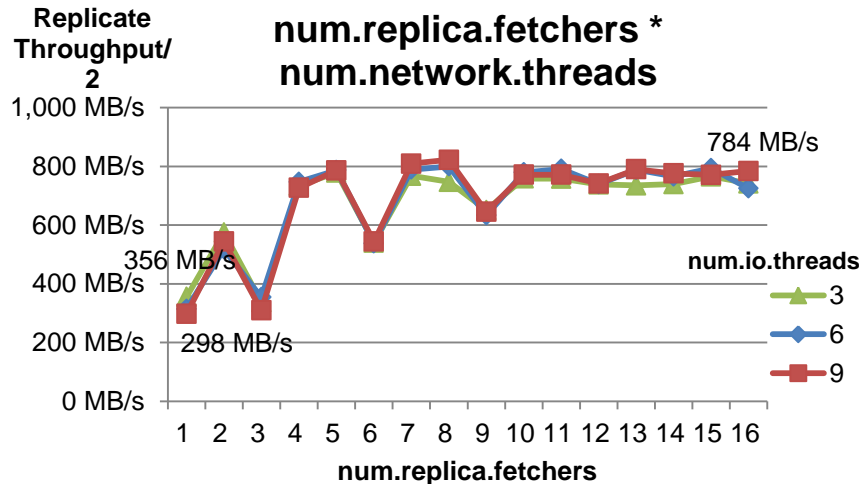
Throughput ② replica.fetch.max.bytes



- ◆ replica.fetch.response.max.bytes=10→100MBでReplicateスループットは 278→287 MB/sに若干向上
- ◆ replica.fetch.max.bytes=1→5MBでReplicateスループットは 287→356 MB/sに向上

③ Broker間のレプリケーション性能のチューニング

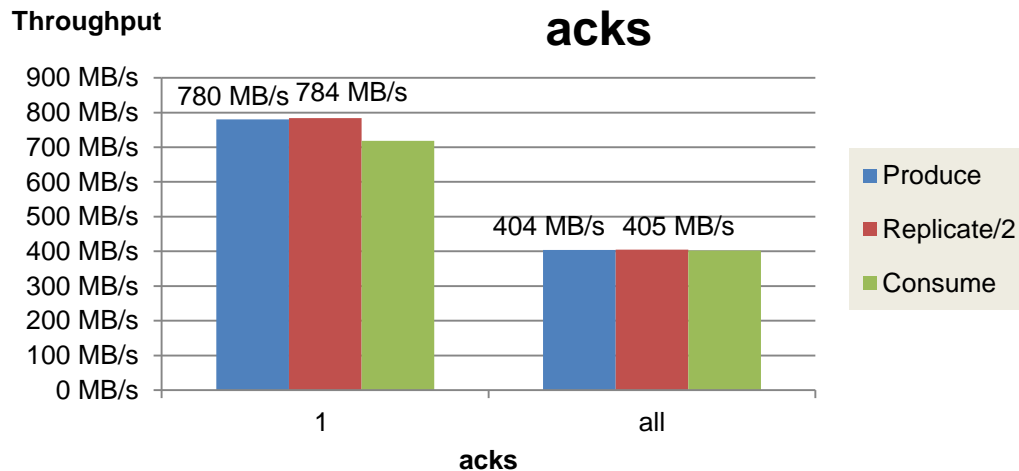
- Replica fetcherのスレッド数と、Socket Serverのスレッド数を同時にチューニング



- ◆ num.replica.fetchersのチューニングは効果は大きいですが、num.io.threadsの影響は少なかった
- ◆ num.replica.fetchersがBrokerの倍数(3の倍数)の時は、一部のReplica Fetcherにすべて自BrokerのPartitionを割り当ててしまうため、そのReplica fetcherが使用されず性能が低下
- ◆ num.io.threads=3→9, num.replica.fetchers=1→16でReplicateスループットは359→784MB/sに向上

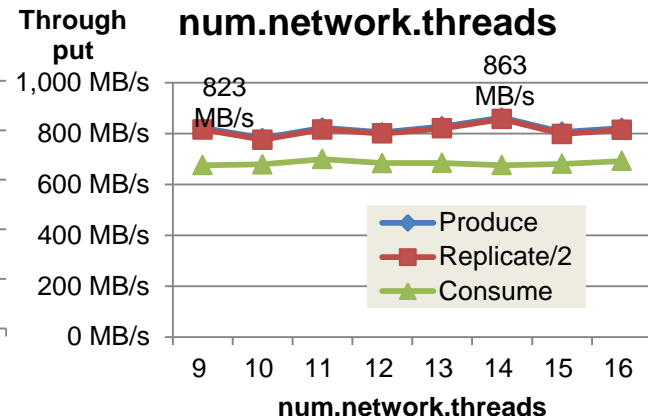
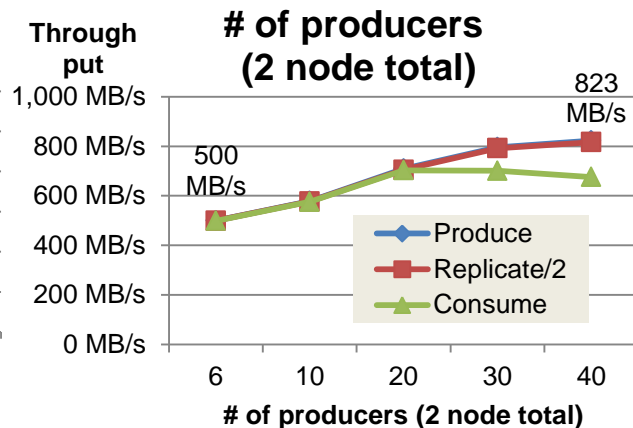
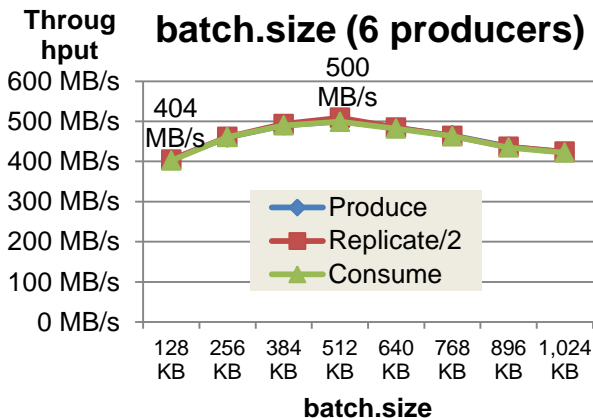
③ Broker間のレプリケーション性能のチューニング

- acks=1 -> all に変更して、レプリケーションを非同期 -> 同期に変更



acks=1の時点でProduceとReplicateのスループットがほぼ釣り合っていたにも関わらず、acks=allに設定するとスループットは大幅に低下してしまった

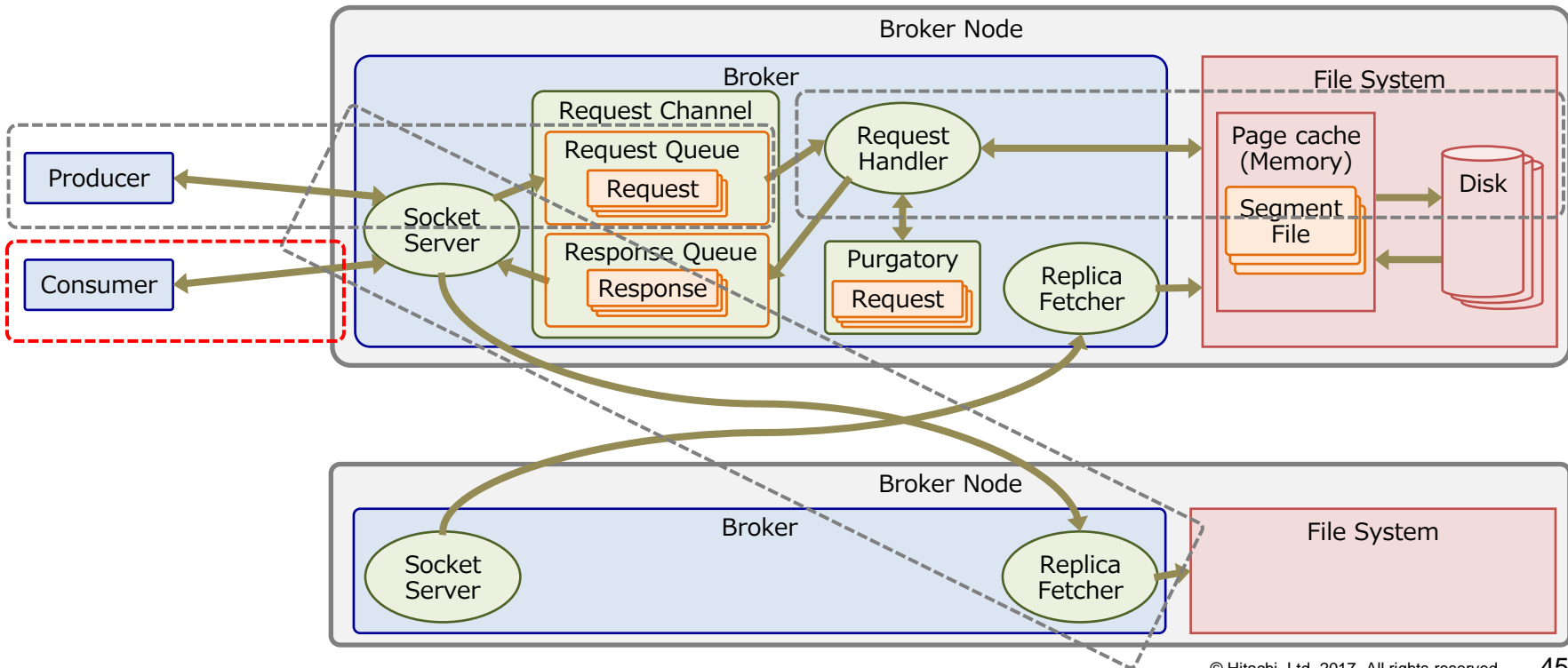
• acks=allに設定してProducer送信性能を再チューニング



- ◆ batch.size= 128->512KBで、Produceスループットは 404->500 MB/s に向上
 - acks=allではRecord Batchサイズhを増やして1回のリクエストサイズを大きくした方が有利？
- ◆ Producer数= 6->40で、Produceスループットは 500->823 MB/s に向上
 - 1Producerあたりユーザ/ネットワークスレッドで2コア使用するため、最大40Producer(80コア)
 - acks=allではProducer数(=Brokerとの接続数)を増やすとよい？
- ◆ num.network.threads= 3->14で、Produceスループットは 823->863 MB/s に若干向上

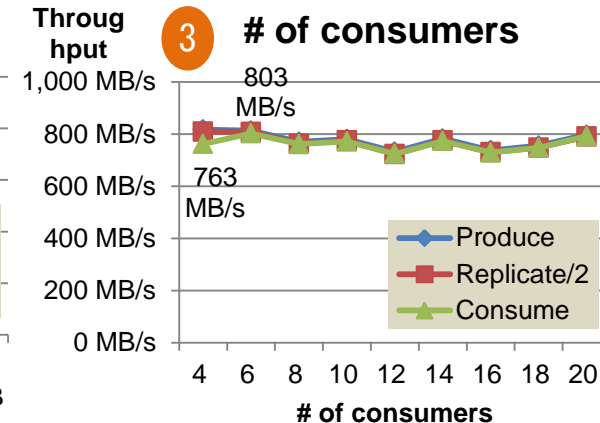
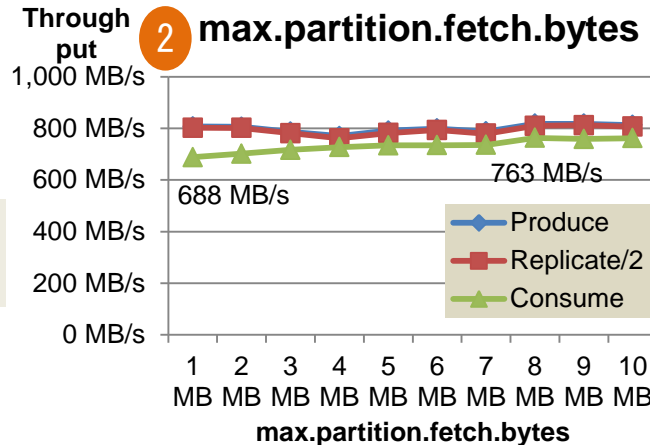
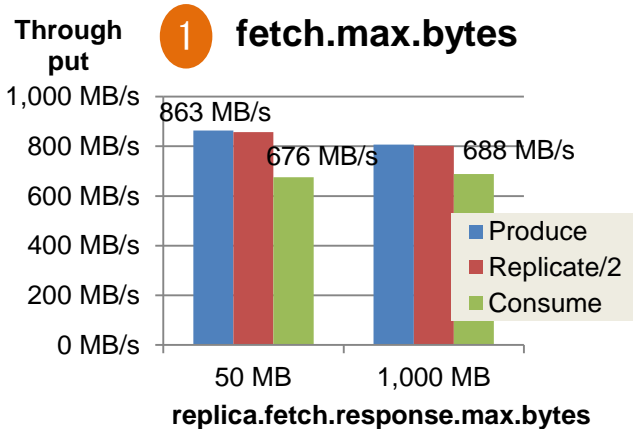
④ Consumerの取得性能のチューニング

- 目的: acks=all で Consumeスループットを最大化する
 - Producerの送信からConsumerの受信まで、システム全体を通したスループットを最適化する



④ Consumerの取得性能のチューニング

- Fetchリクエストの、①合計取得サイズと②Partition単位の最大取得サイズを変更
- ③Consumer GroupのConsumer数を変更



- ◆ fetch.max.bytes= 50→1000MBで、Consumeスループットは 676→688 MB/s に若干向上
 - ◆ max.partition.fetch.bytes= 1→8MBで、Consumeスループットは 688→763 MB/s に向上
 - ◆ Consumer数=4→6で、Consumeスループットは 763→803 MB/s に若干向上
- ⇒ Produce/Replicate/Consumeスループットがほぼ一致し、システム全体で 803 MB/sを達成

Producerに追加したレコードをConsumerで取得するまでの時間を測定

② Producerのバッファリング時間: 50s

- 1Producerの送信レートは $803\text{MB/s} \div 40 \text{ Producer} \doteq 20\text{MB/s}$
- 1Producerのbuffer.memoryは 1024MB
- よって、レコード追加速度 > Producerの送信速度 の場合はバッファが満杯となり、 $1024 \text{ MB} \div 20\text{MB} \doteq 50\text{s}$

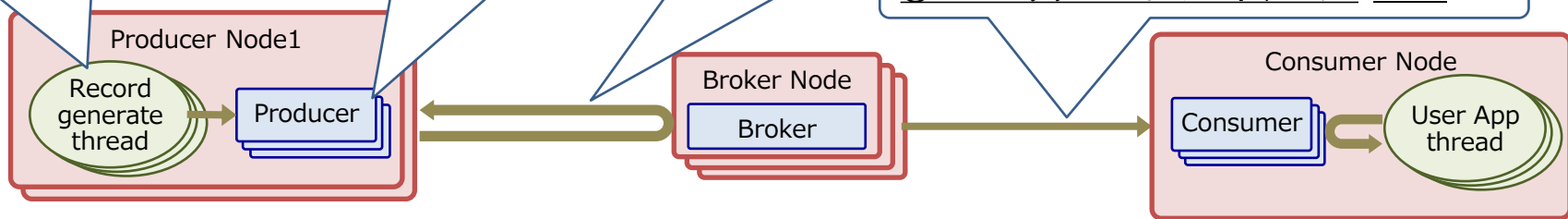
① Send APIのブロック時間: 100ms

Producerのバッファ空き待ち

③ Produceリクエストのレイテンシ: 2.5s

acks=allなのでレプリケーションを含む

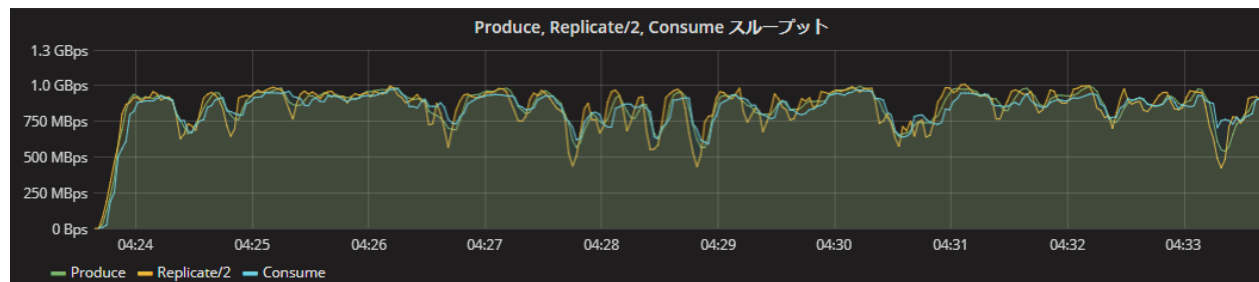
④ Fetchリクエストのレイテンシ: 150ms



- レコード追加速度がProducer送信性能を超えてバッファが詰まるため、全体のレイテンシは 約53秒
- レコード追加速度をKafkaの最大スループット(800MB)以下に抑えれば、レイテンシは 3秒程度 と予想
 - レコード蓄積を待つためにlinger.msを増やす必要があるかもしれない

5. まとめ

- 3レプリカの同期レプリケーション(acks=all)で平均 803MB/sのスループットを達成
 - ネットワーク帯域の理論値(1,170MB/s)の約70%
 - スループットは時間とともに上下しており、1,000MB/s(理論値の85%)程度の帯域は使用している



チューニングのポイント

- Producer
 - batch.size と Producer数(=Brokerとの接続数)を増やす
- Broker
 - num.replica.fetchers を増やす
(ただしReplica fetcherに対するPartition割り当ての偏りに注意)
- Consumer
 - FetchサイズとConsumer GroupのConsumer数を増やす

- Apache Kafka, Apache ZooKeeperは、Apache Software Foundationの米国およびその他の国における登録商標または商標です。
- Javaは、Oracle Corporation及びその子会社、関連会社の米国およびその他の国における登録商標です。
- HITACHIは、株式会社 日立製作所の商標または登録商標です。
- その他記載の会社名、製品名などは、それぞれの会社の商標もしくは登録商標です。