

Advanced
Network
Architecture
Research

次世代高速・高品質インターネット のための技術課題



村田 正幸
大阪大学 大学院基礎工学研究科 情報数理工学専攻
先進ネットワークアーキテクチャ研究室
e-mail: murata@icses.osaka-u.ac.jp
http://www-anaicses.osaka-u.ac.jp/ murata/

Advanced
Network
Architecture
Research

インターネットの高速・高品質化に 向けた技術課題；概観

- ネットワークに対する2つの見方
- 究極のアーキテクチャは存在するか？
- ATMの反省
- インターネットQoS；よくある間違い

Advanced
Network
Architecture
Research

テレコムコミュニティの見方

- 「ネットワークが提供すべきは信頼性のあるコネクション型」
 - X データを送る前にコネクション設定を行い、識別子を端末に与える
 - X コネクションを設定した後、パラメータ、品質、コストに対する交渉を行う
 - X 双方向通信、順番を保証したパケット転送を行う
 - X 輻輳制御機能を提供する
- 主人公はネットワーク
- アプリケーション
 - X 電話、動画（テレビ会議）
 - X リアルタイム（人と人とのコミュニケーション）

M. Murata 3

Advanced
Network
Architecture
Research

インターネットコミュニティの見方

- 「ネットワークの仕事はビットを運ぶこと」
 - X いくらがんばってもネットワークが信頼性を確保することは難しい
 - X ホストはそれを受け入れてエラー制御をおこなう
 - X フロー制御は自分でする
- 主人公はコンピュータ
- アプリケーション
 - X telnet, ftp, WWW (http)
 - X パースト的、リクエスト/レスポンス

M. Murata 4

Advanced
Network
Architecture
Research

複雑な制御をどこにおくか？

- コネクションレス型 トランスポート層（ホスト）
 - X ホストの処理能力の向上
 - X 信頼性より高速な転送が必要なアプリケーションもある
- コネクション型 ネットワーク層（ノード）
 - X ユーザが希望するのは信頼性の高いトラブルのないサービス
 - X 実時間音声や動画はコネクション型のほうが簡単

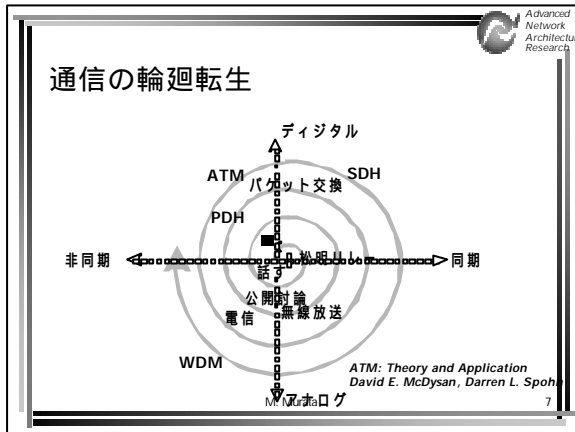
M. Murata 5

Advanced
Network
Architecture
Research

究極の通信技術は存在するか？

- 「ATMはマルチメディア情報を統一的に扱う究極の統合通信網である」
 - X マルチメディア情報とは何か？
 - X 統合通信網は必要か？
 - X エンドシステムは「端末」か？
- ネットワークの目的
 - X 情報共有の手段（含：マルチメディア情報）
 - ✓ 情報発信、情報共有、情報流通、新しい（仮想）コミュニティの形成
 - X コミュニケーションの手段
 - ✓ 音声から動画へ、メディアの多様化

M. Murata 6



- ### ATM (=Telecom)の反省
- 統計多重効果をはやくあきらめるべきだった
 - 「何もかも扱う」という命題を背負った（背負わされた？）ところに問題がある
 - データ通信には原理的に向かない（次ページ参照）
 - X データ端末はパケット交換か？
 - エンドシステムの軽視
 - X コンピュータは端末ではない
 - ATMのAPIは解放されていたか？
 - X インターネットが目の前にあったからこそ、それに通じたWebというアプリケーションが生まれた
 - ✓ 背景：画像圧縮技術、GUI、画像表示能力
 - ✓ にわとりと卵（？）
- M. Murata



- ### アプリケーションは重要か？
- ネットワークインフラにとってアプリケーションは
 - X ショーウィンドウ
 - X 画像アプリケーションが回線を埋め尽くすわけではない
 - 重要なのはサービスアーキテクチャ
 - X エンドシステム・アクセス系・バックボーンのどこで実現するかが問題
 - 一億総プログラマの時代
 - X 趣味・嗜好の細分化 ネットワーク提供者がアプリケーションを考えるのは所詮無理
 - X ユーザが自由にアプリケーションを作れるような環境を提供することが重要
 - X 何が登場するかわからない 単純なネットワーク構造が必要
- M. Murata

- ### データ系QoS：よくある間違い
- データ系のQoSはパケット棄却率・遅延
 - データ系のQoS保証項目はアプリケーションレベルでの遅延
 - アーラン呼損式に相当するのは待ち行列（網）理論に基づく結果
 - TCPがあればネットワーク内部に輻輳制御メカニズムはいらない
 - TCPはプロトコルとして同じ振る舞いが規定されているので、公平な通信サービスが提供される
 - TCPはもはや古いので、新しい軽量・高性能プロトコルが必要
 - WDMを導入すれば、ネットワーク速度は波長数分向上する
- M. Murata

- ### 高速・高品質化に向けた技術課題
- 課題1：実時間系QoS保証
 - 課題2：バックボーンの高速度化
 - 課題3：プロトコルの高速度化
 - 課題4：エンドホストの高速度化
 - 課題5：ネットワーク機能の再配分
 - 課題6：公平性の問題
 - 課題7：ネットワークプロビジョニング
 - 課題8：基礎理論はあるか？

Advanced Network Architecture Research

データ系QoSとは？

- ユーザの我慢の上に成り立っている
- データ系の通信保証は原理的に不可能
 - X 「64Kbpsを保証する」 = 「アクセス回線のみ保証」 or 「呼損の発生」 or 「64Kbps以上は許さない」
- データ系サービスのQoS
 - X 少なくともパケット遅延・棄却率ではない
 - X ユーザレベルの遅延を「高速化」する
例：Webのドキュメントダウンロード遅延
- 「インターネットサービス」を受け入れるには世代交代が必要(？)

M. Murata 31

Advanced Network Architecture Research

課題2：バックボーン的高速化

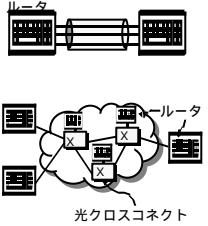
- WDMネットワークにすべてのネットワーク機能を取り込むのには無理がある
 - X 例：コネクション設定、輻輳制御、経路制御
- ただし、IPの機能を一部取り込むことによって高信頼化・高性能化を実現することは可能
 - X 高信頼化：代替経路を設定
 - ✓ link protection
 - ✓ dedicated-path protection
 - ✓ shared-path protection
 - X 高性能化：ルータによるパケットフォワード直接パスの設定
 - X 前提：ネットワークディメンジョンングがうまく機能すること
 - ✓ 波長ルーティングによる容量変更

M. Murata 32

Advanced Network Architecture Research

フォトニックインターネットアーキテクチャ

- 4つのフォトニックネットワークアーキテクチャ
 1. WDMリンクネットワーク
 - ✓ ルータ間をWDMリンクで接続
 2. WDMバスネットワーク
 - ✓ ルータ接続(フォトニックバスによる論理ネットワークの構築)

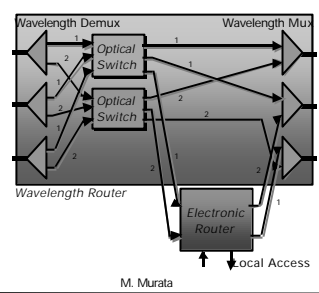


光クロスコネクタ

M. Murata 33

Advanced Network Architecture Research

波長ルータの構成

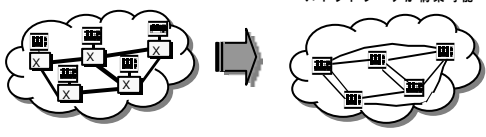


M. Murata 34

Advanced Network Architecture Research

波長バスによる論理トポロジーの構成

- 物理トポロジー
- 論理トポロジー
 - ✓ ルータから見ると冗長度の高いネットワークが構築可能



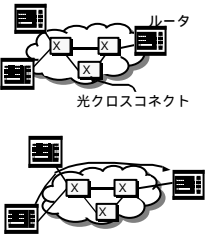
- 機能分担(信頼性技術)をどうするか?
 - X IP層；経路制御
 - X WDM層；バスプロテクション

M. Murata 35

Advanced Network Architecture Research

フォトニックインターネットアーキテクチャ(続)

- 4つのフォトニックネットワークアーキテクチャ
 3. WDMバスネットワーク
 - ✓ ラベルスイッチング(波長をラベルと見たMPLS)
 4. WDMパケットネットワーク
 - ✓ パーススイッチング(オンデマンド波長・経路選択)




光クロスコネクタ

M. Murata 36

Advanced Network Architecture Research

1000波長WDMの実現可能性とその効果

- 与条件
 - X 物理トポロジー
 - X トラフィック量、トラフィック分布
 - 基本：電話網
 - スケールファクタにより調整
- 物理トポロジーから論理トポロジーを生成
 - 必要波長数、ルータのバケット処理能力
- 村田正幸、北山研一、宮原秀夫（大阪大学）「1000波長WDMによるインターネットにおけるネットワークボトルネックの解消」発表予定



<http://www.nttco.jp/databook/setubi/>

M. Murata 37

Advanced Network Architecture Research

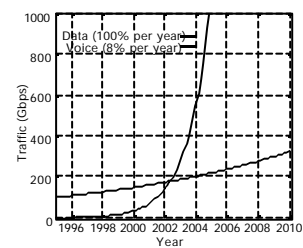
トラフィックマトリックス

- NTTの電話網を基本とする
- X NTT情報Webステーション
 - <http://www.ntt-east.co.jp/info-st/network/traffic/index.html>
- X 特徴
 - 近隣県間のトラフィック量 遠距離県間のトラフィック量
 - Zipfの法則(?)
 - 東京への一極集中
 - 大阪 兵庫：大阪 東京：大阪 岩手=300:100:1
 - トータル 30Gbps (1アーランを64Kbpsで換算)
 - スケールファクタを導入
- X データ系トラフィック?

M. Murata 38

Advanced Network Architecture Research

インターネットトラフィックの成長

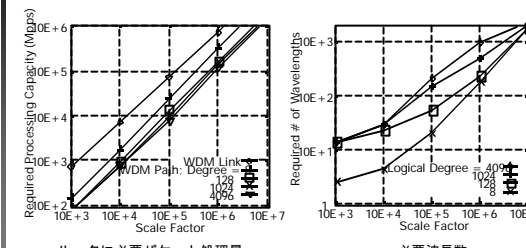


- USの例
 - X 音声：年8%
 - X データ：年100%
 - 3年で1桁上昇
 - X K.G. Coffman and A.M. Odlyzko, "The size and growth rate of the Internet," <http://www.research.att.com/~amo>

M. Murata 39

Advanced Network Architecture Research

1000波長WDMによる効果



ルータに必要なバケット処理量

必要波長数

□ ルータボトルネック!

M. Murata 40

Advanced Network Architecture Research

課題3：プロトコルの高速化

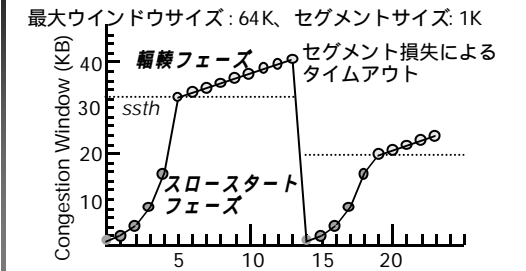
- 従来の議論
 - X プロトコル処理のハードウェア化、パラレル化
 - X 軽量プロトコルの実現
 - X エラー制御と輻輳制御の分離
 - Proprietaryなプロトコルはもはや通用しない
- TCPの性能向上?
 - X TCP Tahoe TCP Reno TCP NewReno TCP SACK TCP Vegas (?)
 - X プロトコルマイグレーションの確保
 - 送信側の変更は可能
 - バージョンの異なるTCP混在時の性能?

M. Murata 41

Advanced Network Architecture Research

TCPのウィンドウサイズの動き

最大ウィンドウサイズ：64K、セグメントサイズ：1K



Congestion Window (KB)

時刻 (単位: RTT)

転換フェーズ

セグメント損失によるタイムアウト

sssth

スロースタートフェーズ

M. Murata 42

Advanced Network Architecture Research

TCPのバージョン

- たまたま1個だけセグメントを落としただけならセグメントをすぐに再送
fast retransmit; TCP Tahoe
- ウィンドウサイズを半分にする
fast recovery; TCP Reno
- Go-Back-NからSelective Repeatへ; TCP SACK
- RTTに基づいてウィンドウサイズを調整する: TCP Vegas

$$diff = cwnd(t) / basertt - cwnd(t) / observed_rtt$$

$$cwnd(t+1) = \begin{cases} cwnd(t) + 1, & \text{if } diff < a / base_rtt \\ cwnd(t), & \text{if } a / base_rtt \leq diff < b / base_rtt \\ cwnd(t) - 1, & \text{if } b / base_rtt \leq diff \end{cases}$$

M. Murata 43

Advanced Network Architecture Research

TCP VegasとTCP Reno混在時のスループット比較

- TCP Vegasだけなら ~ 40%のスループット向上

Throughput[Kbps]

Buffer Size[packets]

TCP Reno: 5本
10Mbps
1.5Mbps
TCP Vegas: 5本

M. Murata 44

Advanced Network Architecture Research

課題4: エンドシステムを含めた性能?

- パケット伝送時間がネットワーク性能を決めるわけではない
- X 回線容量
- X ルータ処理能力
- X ホスト(サーバ・クライアント)処理能力
- X アプリケーション

M. Murata 45

Advanced Network Architecture Research

アプリケーションの遅延配分

Webページのダウンロードの場合

- エンドシステムの重要性
- ネットワーク高速化の限界
- 重要なのはバランスのとれた資源配分

ネットワーク転送 37%

DNS 15%

コネクション設定 28%

サーバ処理 20%

参考文献: 藤田 靖征, 村田 正幸, 宮原 秀夫, "Webサーバシステムのモデル化と性能評価," 電子情報通信学会論文誌, 1998.

Produced from <http://www.telcordiacom/pub/huitema/stats>

M. Murata 46

Advanced Network Architecture Research

高速プロトコル処理システム

- ソケット層プロトコルの高速化 (バッファ管理)
- エンドシステムのプロトコル処理高速化 (ゼロコピー)
- ネットワークにおけるプロトコル高速化 (TCPの輻輳制御メカニズム)
- ネットワークの高速化 (IP over WDM)

Socket Layer

TCP Layer

IP Layer

Data Link/Physical Layer

2.4Gbps

MAPOS: 20Gbps

M. Murata 47

Advanced Network Architecture Research

課題5: 通信処理機能の再配分

- エンドホストに頼りすぎ
 - X TCPによる輻輳制御
 - ✓ 輻輳制御は本来ネットワーク機能
 - X 「公平なサービス」を阻害する
 - ✓ 輻輳制御を行わないホスト(バグ、コードの改変)
 - ✓ 「サービスの有料化」に対する障害
- 通信処理機能の再配分
 - X フロー制御、誤り制御、輻輳制御、経路制御
 - X RED、DRR、ECN、diff-serv、int-serv (RSVP)、ポリシールーティングは輻輳制御のネットワークへの回帰
 - X ただし、過度の回帰はインターネットのメリットをなくす

M. Murata 48

課題6：公平性の問題

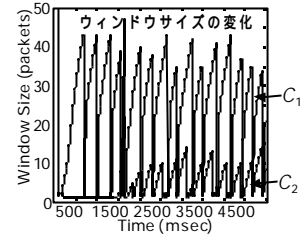
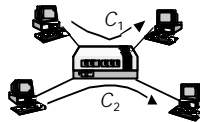
- 資源が無限大になる時代がやってくることはない！
 - × 帯域を埋めるようなアプリケーションがあれば、ネットワークが機能しないのは自明
 - × データ系は帯域に合わせて送信する
- 帯域をいかに公平に分配するか？
 - × TCPの輻輳制御は公平ではない
 - ✓ いったんウィンドウサイズを下げると（短期的に）大きくならない
 - ✓ RTT、帯域による不公平性
 - ➡ ルータによる処理の必要性；RED、DRR
 - × 実時間アプリケーションとデータ系アプリケーションの帯域の配分？

M. Murata

49

TCPコネクション間の公平性： 伝播遅延時間が異なる場合

- $2\tau_1 = \tau_2$ の場合

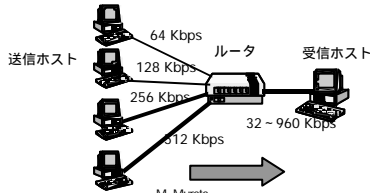


M. Murata

50

TCPコネクション間の公平性： 回線容量が異なる場合

- 出力回線容量に対して、接続回線容量に応じたスループットが得られるか？
 - × Relative Throughput = 回線容量を1とした場合のスループット

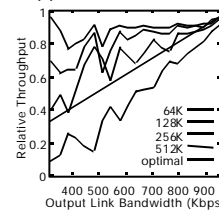


M. Murata

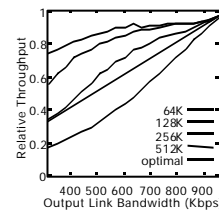
51

TCPコネクション間の公平性： 回線容量が異なる場合の結果例

- ルータがDrop Tailの場合
- ルータがREDの場合



M. Murata



52

課題7：ネットワークプロビジョニング

- 現状はネットワークトラフィックの振る舞いの把握と分析
 - × "Cooperative Association for Internet Data Analysis," <http://www.caida.org/>
 - × "Internet Performance Measurement and Analysis Project," <http://www.merit.edu/ipma/>
- 意味のあるデータを拾い出せるか？
 - × 経路制御による非安定性
 - × TCPのセグメント再送
 - × 実時間アプリケーションのレート適応・遅延適応制御
 - × 低利用率は
 - ✓ 輻輳制御のため？
 - ✓ 低速アクセス回線のため？
 - ✓ 低速エンドホストのため？

M. Murata

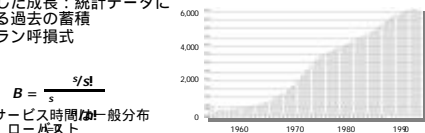
53

アーラン呼損式はなぜ信用できるのか

目標呼損率 回線容量
トラフィック測定 回線増強
呼損率 = ユーザ品質
安定した成長：統計データに関する過去の蓄積
アーラン呼損式

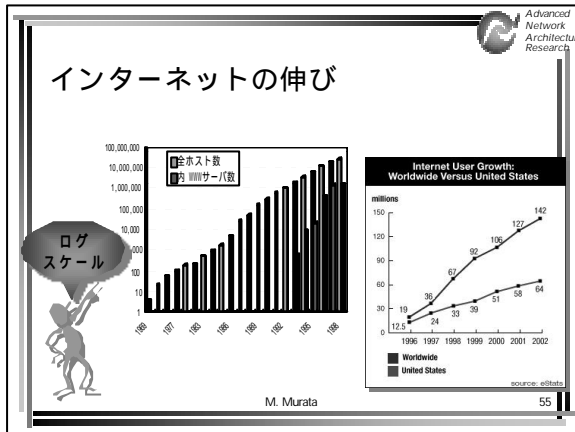
$$B = \frac{y/s}{1-y/s}$$

- i. サービス時間一般分布ローバースト
- ii. ボアソン到着は普遍的な事実
品質測定 = 呼損率を測る



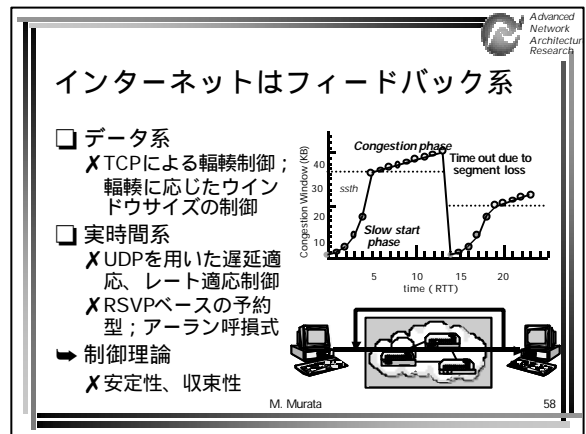
M. Murata

54



- ### ネットワークプロビジョニングの課題
- (少なくとも) ネットワークプロビジョニングレベルでの品質予測
 - 回線交換にない新たな問題
 - X 品質とは何か?
 - X 品質を測定できるか?
 - X サービス品質の課金への反映?
 - X マルチメディアトラフィックの予測の困難性
 - X ネットワーク測定では不十分
 - ネットワークトラフィック測定 分析 回線容量設計
 - X フィードバックループを前提としたネットワーク設計論の確立
 - X 柔軟な帯域設定を持つネットワークが大前提 (ATM、フォトニックネットワーク)
 - X Webトラフィック、Webサーバ、ネットワーク遅延 (RTT) のモデル化
- M. Murata 56

- ### 課題8: パケット交換網の基礎理論?
- M/M/1 (待ち行列網理論) パラダイムは役に立つのか?
 - X わかるのはノードにおけるパケット待ち時間、棄却率
 - ➔ しかし、データ系のQoSはノードにおけるパケット待ち時間ではない
 - X アーラン呼損式 (= 電話網) では 呼損率 = ユーザレベルQoS
 - ルータでの振る舞い?
 - X フィードバック系が上位レベルにある時の振る舞いが重要
 - ユーザレベルのQoS?
 - X アプリケーションレベルの性能指標が重要
 - ✓ e.g. Webドキュメント応答時間
-
- M. Murata 57



- ### ネットワークの基礎理論の構築 (1)
- 輻輳制御
 - X 時間に依存したふるまい
 - X フィードバックシステム: 制御理論
 - フィードバック系を前提にしたトラフィック測定・分析・ネットワークディメンジョンング ネットワーク回線容量設計
 - エンドシステムを包含したQoSアーキテクチャの構築
- ↓
- 回線容量設計をするには
 - ➔ いかにアプリケーション、エンドシステムの影響を除いた測定、モデル化を実現するか?
 - 高品質なネットワークを構築するには
 - ➔ アプリケーション、エンドシステムを含めたQoS設計
- M. Murata 59



ネットワークの基礎理論の構築 (3)

□ 複雑化・巨大化したシステムの評価手法

- X 理論、シミュレーション
 - X 待時系・即時系混合待ち行列網
 - ✓ 即時系コネクションは各ノードの資源を同時に使用
 - ✓ 残りを待時系パケットが使用
 - X 評価モデル？
 - ✓ シミュレーション、解析
- 例：バックボーンにおけるTCPのふるまい？