

Scatter

Definition

- One process sends different messages to each other process
 - Data flow is identical to that of MPI_Gather

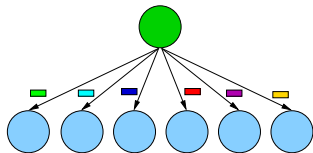
Scatter

Definition

- One process sends different messages to each other process
 - Data flow is identical to that of MPI_Gather

Strategies

- Flat Trees - default
- Binomial Trees
- Chains



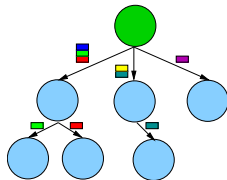
Scatter

Definition

- One process sends different messages to each other process
 - Data flow is identical to that of MPI_Gather

Strategies

- Flat Trees - default
- **Binomial Trees**
- Chains



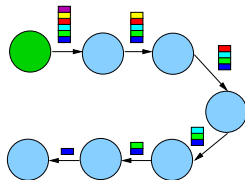
Scatter

Definition

- One process sends different messages to each other process
 - Data flow is identical to that of MPI_Gather

Strategies

- Flat Trees - default
- Binomial Trees
- **Chains**



Modelling the Scatter Operation

Modelling Scatter

- Flat Tree - every process is directly contacted by the *root*
- Other strategies - messages are relayed through auxiliary nodes

Implementation Technique	Communication model
Flat Tree	$(P - 1) \times g(m) + L$
Chain	$\sum_{j=1}^{P-1} g(j \times m) + (P - 1) \times L$
Binomial Tree	$\leq \lceil \log_2 P \rceil \times L + \sum_{j=0}^{\lceil \log_2 P \rceil - 1} g(2^j \times m)$ $\approx \lceil \log_2 P \rceil \times L + (P - 1) \times g(m)$

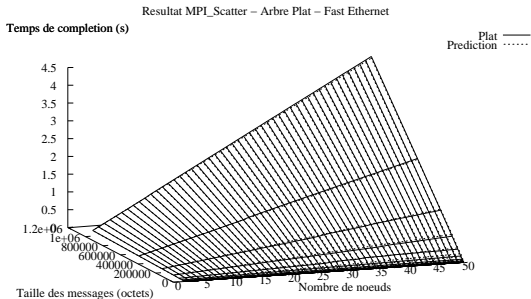
Modelling the Scatter Operation

Modelling Scatter

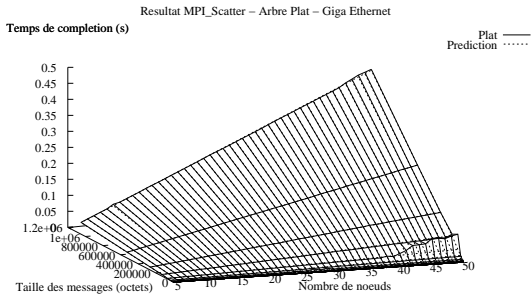
- Flat Tree - every process is directly contacted by the *root*
- Other strategies - messages are relayed through auxiliary nodes

Implementation Technique	Communication model
Flat Tree	$(P - 1) \times g(m) + L$
Chain	$\sum_{j=1}^{P-1} g(j \times m) + (P - 1) \times L$
Binomial Tree	$\leq \lceil \log_2 P \rceil \times L + \sum_{j=0}^{\lceil \log_2 P \rceil - 1} g(2^j \times m)$ $\approx \lceil \log_2 P \rceil \times L + (P - 1) \times g(m)$

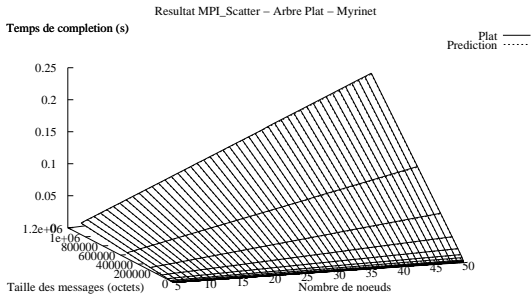
Validating the model - Fast Ethernet



Validating the model - Giga Ethernet



Validating the model - Myrinet



Scatter for Grid Environments

Problem

- Optimisation is harder than Broadcast
- Multi-layered scheduling implies extra efforts

Scatter for Grid Environments

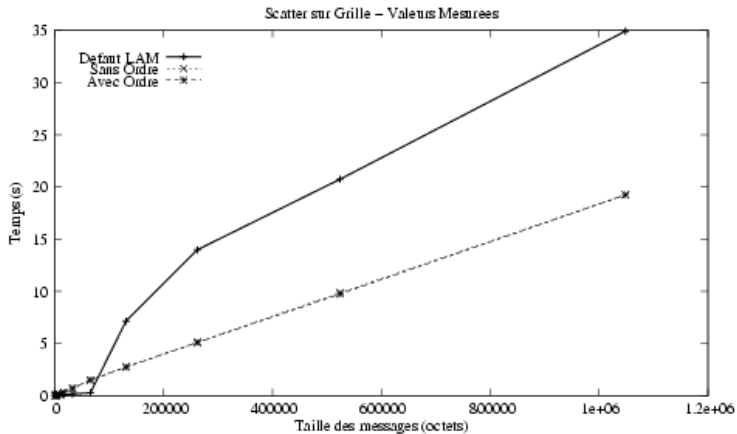
Problem

- Optimisation is harder than Broadcast
- Multi-layered scheduling implies extra efforts

Simple Strategies

- Direct Connexion - standard MPI approach
- Flat Trees with message packing (two layers)
- Flat Trees with order scheduling

Comparison



Summary

Efficient with large messages

- minimises the startup cost on a grid connexion

Future works

- evaluate the impact of memory swap
- improve local area data multiplication (parallel senders)

All-to-All

Personalised Many-to-Many communication pattern

- Every process send a different message to each other

All-to-All

Personalised Many-to-Many communication pattern

- Every process send a different message to each other

Implementation Strategies

- By default, each message is directly sent to the destinations

All-to-All

Personalised Many-to-Many communication pattern

- Every process send a different message to each other

Implementation Strategies

- By default, each message is directly sent to the destinations
 - may saturate the network

All-to-All

Personalised Many-to-Many communication pattern

- Every process send a different message to each other

Implementation Strategies

- By default, each message is directly sent to the destinations
 - may saturate the network
 - contention effects are difficult to model

Modelling All-to-All

Traditional Approach

- Most authors use a model based on the Scatter operation
- Hypothesis: one All-to-All is equivalent to simultaneous Scatters

Modelling All-to-All

Traditional Approach

- Most authors use a model based on the Scatter operation
- Hypothesis: one All-to-All is equivalent to simultaneous Scatters

Drawbacks

- does not suppose network contention
- all processes start communication simultaneously

Our Approach

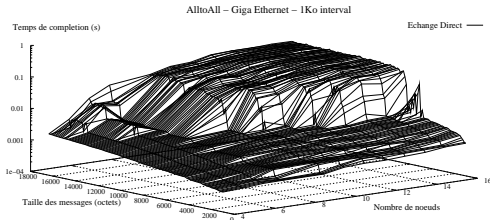
Augmenting the communication model

- Hypothesis: Contention is a linear factor that characterises each network
 - Objective: identify a *Contention Signature* - γ
 - Augment the Scatter model with this signature

$$\begin{aligned} T &= \textit{Free} \times \gamma \\ &= ((P - 1) \times g(m) + L) \times \gamma \end{aligned}$$

Another Factor

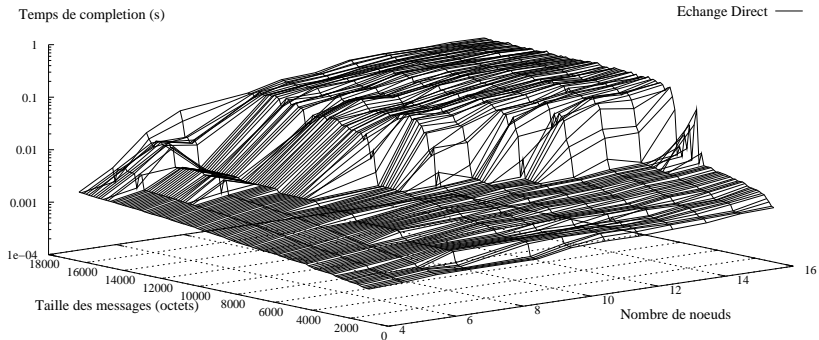
- Buffering, segmentation and memory constraints may impact on the performance
 - Overhead factor δ_P depending on P and m



$$T = ((P - 1) \times g(m) + L) \times \gamma + (P - 1) \times \delta_P$$

Buffering

AlltoAll – Giga Ethernet – 1Ko interval



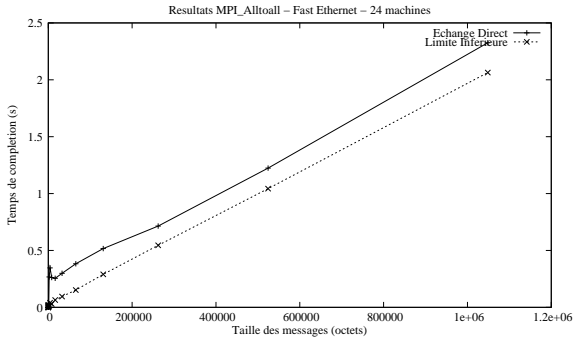
Defining Parameters

To obtain γ et δ_P

- Performance prediction (Scatter model)
- Execution sample
- Parameter fitting with Linear Regression
 - Least Squares method
- Once identified, γ et δ_P are used to correct the predictions

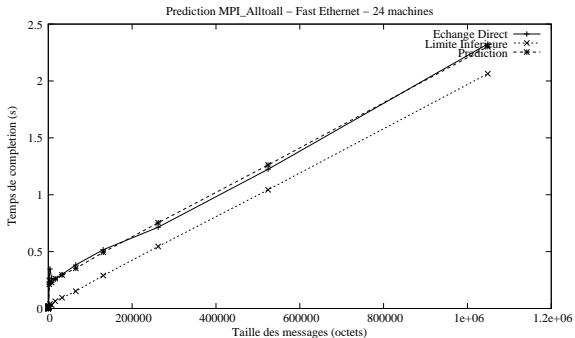
Fast Ethernet Network

- Reduced contention

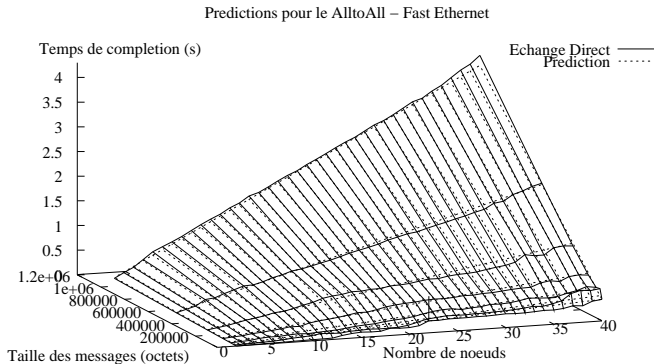


Fast Ethernet Network

- $\gamma = 1.0195$
- $\delta = 8.23$ ms (for messages larger than 2 KB)

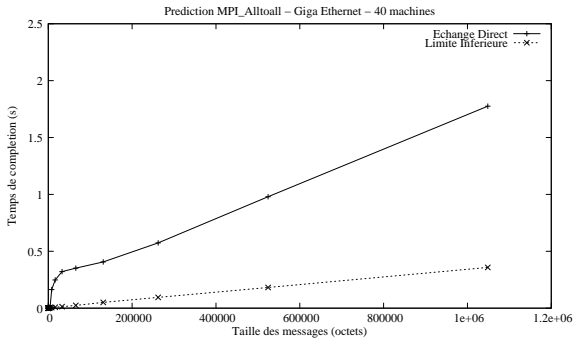


Applying the Model



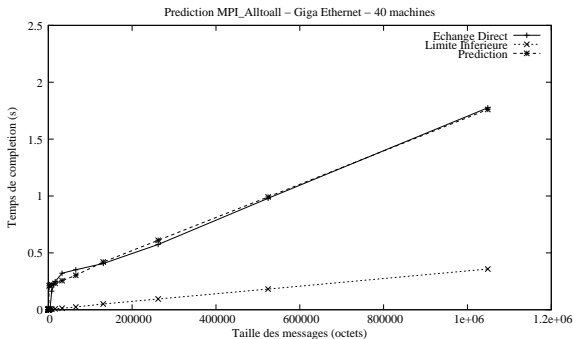
Giga Ethernet Network

- Impact of contention is important
- Buffering plays an important role

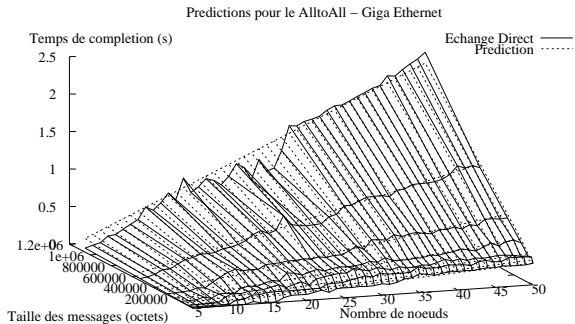


Giga Ethernet Network

- $\gamma = 4.3628$
- $\delta = 4.93$ ms (for messages larger than 8KB)

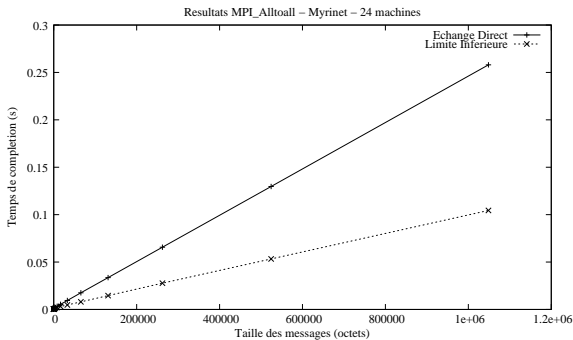


Applying the Model



Myrinet Network

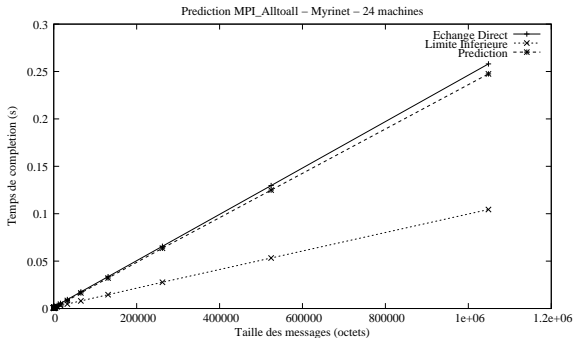
- Contention is important
- No buffering side-effects



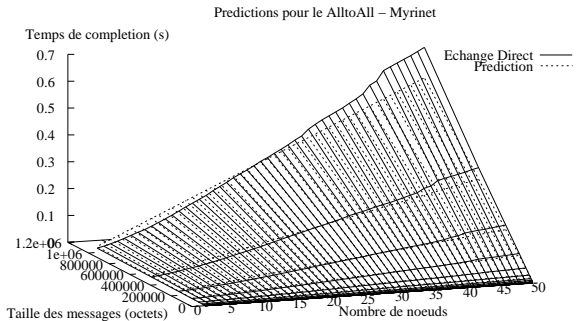
Myrinet Network

- $\gamma = 2.4975$

- $\delta = 0$



Applying the model



AlltoAll

THE hardest problem

- How to explore heterogeneity?

AlltoAll

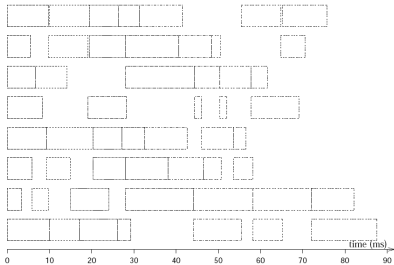
THE hardest problem

- How to explore heterogeneity?

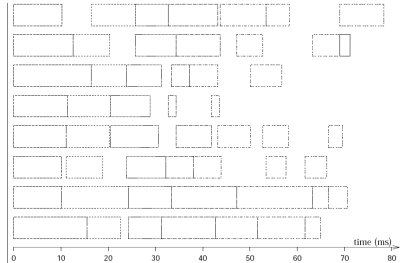
Works on the AlltoAll-v operation

- Messages with different sizes
 - different sizes = different communication times
- Matching heuristics

Motivation



Fixed Pattern



Max-min

Heuristics for the AlltoAll

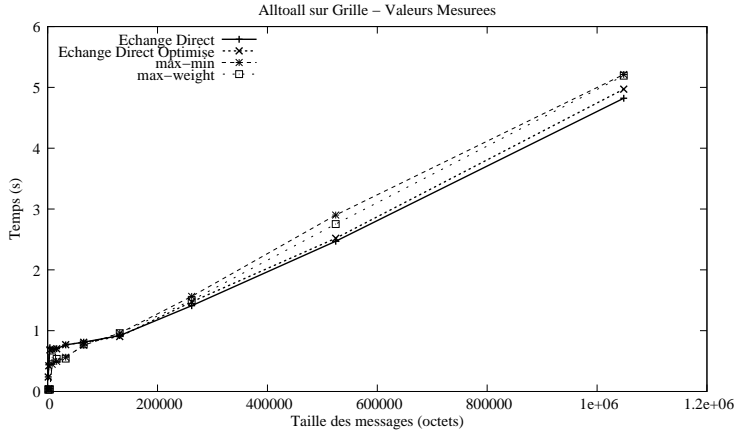
Max-min heuristic

- The minimum communication weight is maximised

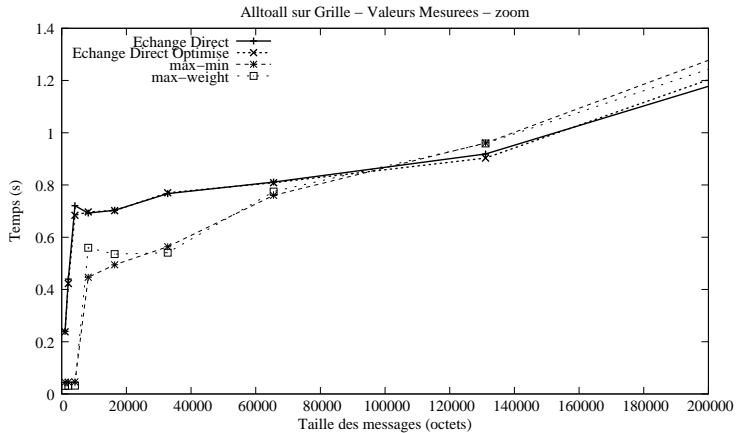
Max-weight heuristic

- The sum of the communication weights is maximised
 - Maximises the data transfer in a round

Comparison



Comparison



Summary

- Limited efficiency
 - Contention may prevents any performance improvement
- Suited for small messages

Future Works

- Extend the study with the *All-to-All-v* operation

Available Bandwidth on Wide-Area Connections

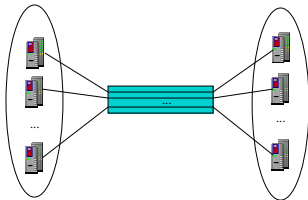
In collaboration with Carlos Barrios (June 2005)

- Modelling massive data transfer rates
- Identify available bandwidth on broadband connexions
- Help to evaluate different tuning parameters

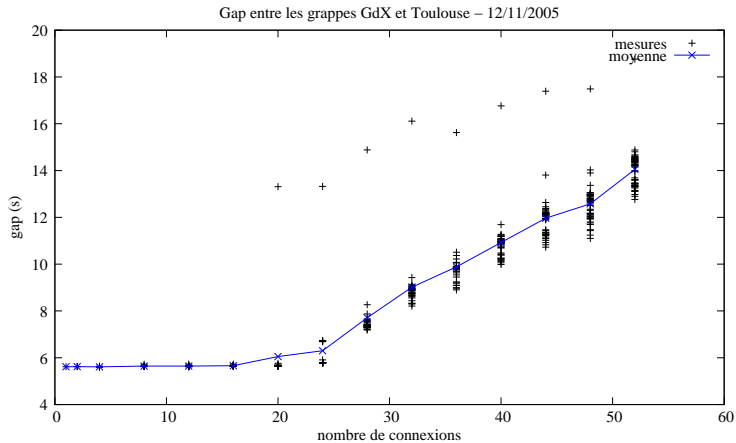
How to saturate a broadband link

Problem: a single connexion is not enough

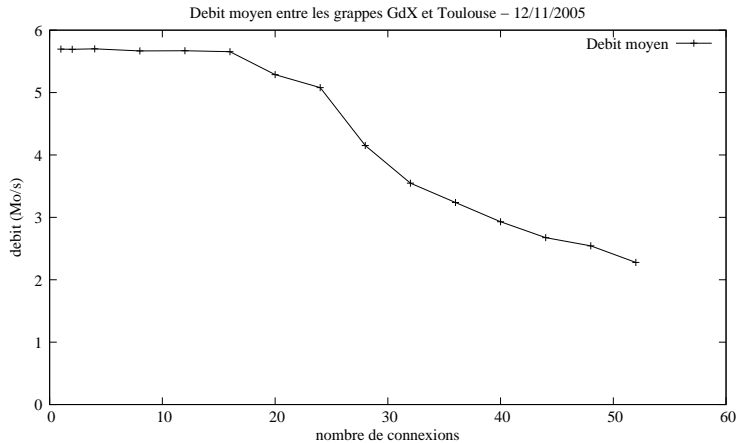
- Simultaneous connexion of several nodes
 - TCP fairness



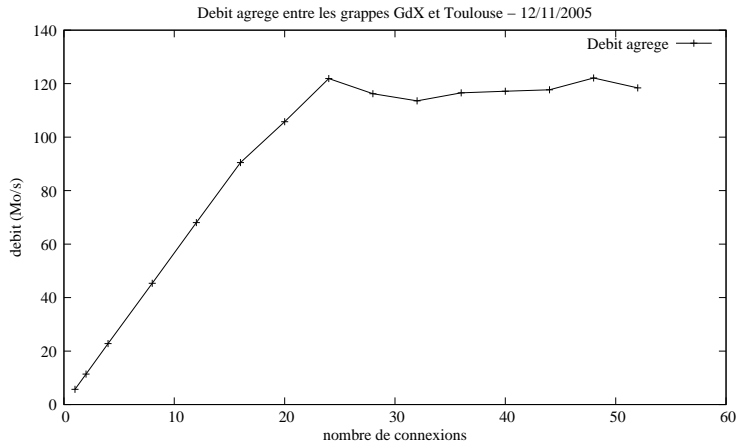
Gap



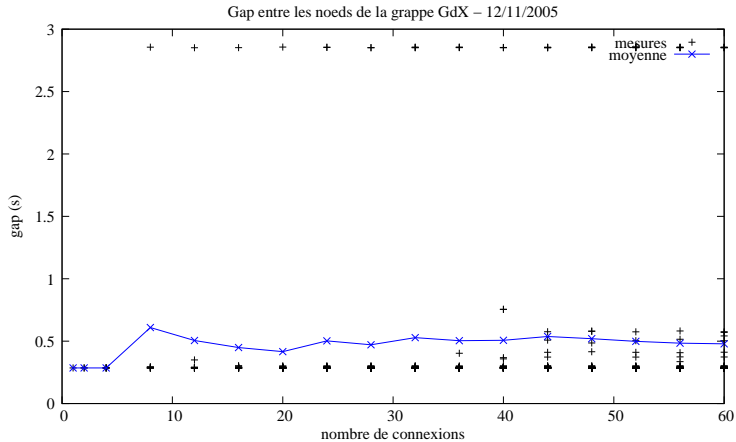
Bandwidth



Aggregate Bandwidth



Quantifying contention on LANs



Bandwidth

