LUDWIG-
MAXIMILIANS-
UNIVERSITÄT
MÜNCHEN

IBE

MEDIZINISCHE FAKULTÄT
INSTITUT FÜR MEDIZINISCHE INFORMATIONSVERARBEITUNG
BIOMETRIE UND EPIDEMIOLOGIE

# Parallelized preprocessing algorithms for
# high-density oligonucleotide array data

Markus Schmidberger
Ulrich Mansmann

HiCOMB 2008
Seventh IEEE International Workshop on High Performance Computational Biology
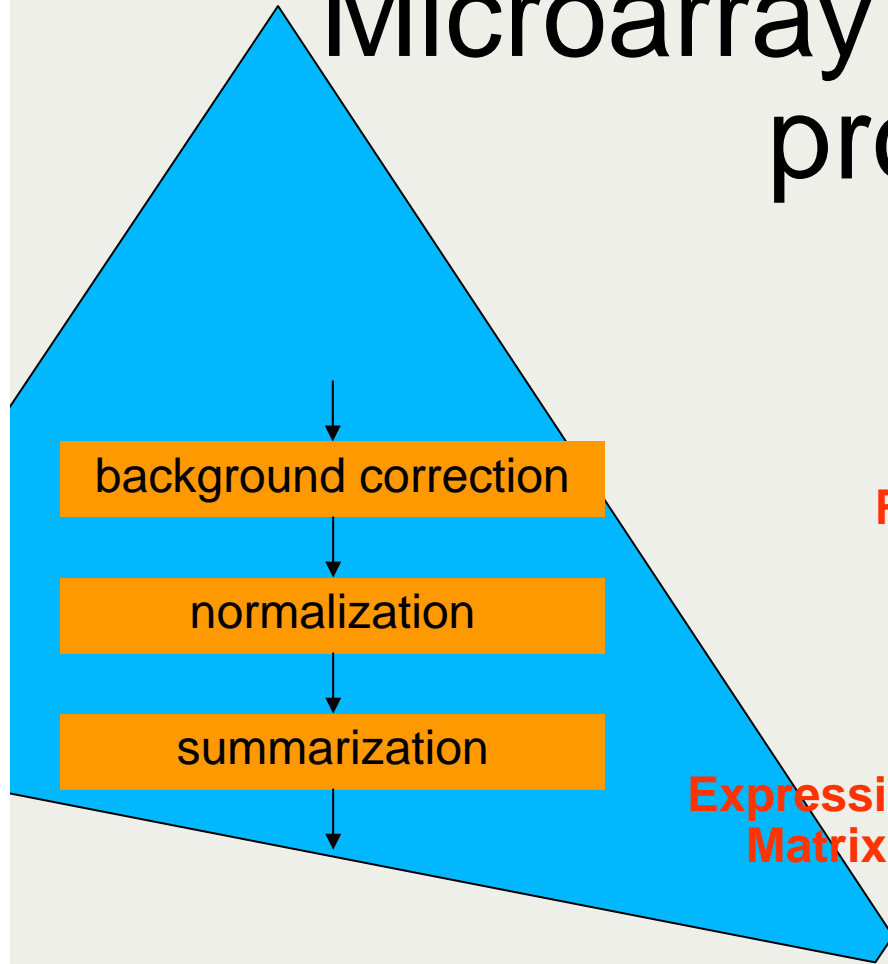April 14, 2008; Miami, Florida, USA

IBE
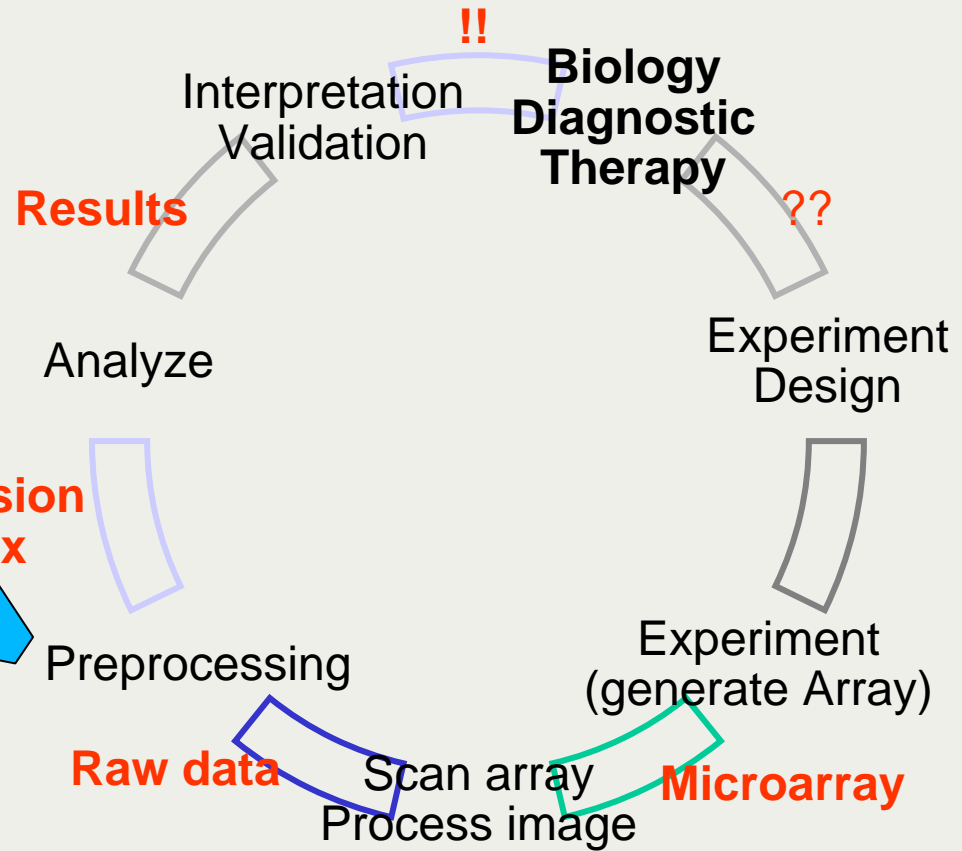http://ibe.web.med.uni-muenchen.de

# Overview

1. Introduction to Microarrays

2. Problems & Challenges & Solutions

3. Parallelization and R

4. Parallelization for preprocessing

5. Results

# Microarray Data Analysis process



background correction

normalization

summarization

**Expression Matrix**

**!!**

Interpretation Validation

**Biology Diagnostic Therapy**

**Results**

**??**

Analyze

Experiment Design

Preprocessing

Experiment (generate Array)

**Raw data**

Scan array Process image

**Microarray**

3

# Sources of errors

amount of RNA in the biopsy
efficiencies of
- RNA extraction
- reverse transcription
- Labeling
- fluorescent detection

- probe purity and length
- distribution
- spotting efficiency, spot size
- cross-/unspecific hybridization
- stray signal

**Systematic**
- similar effect on many measurements
- corrections can be estimated from data

**Stochastic**
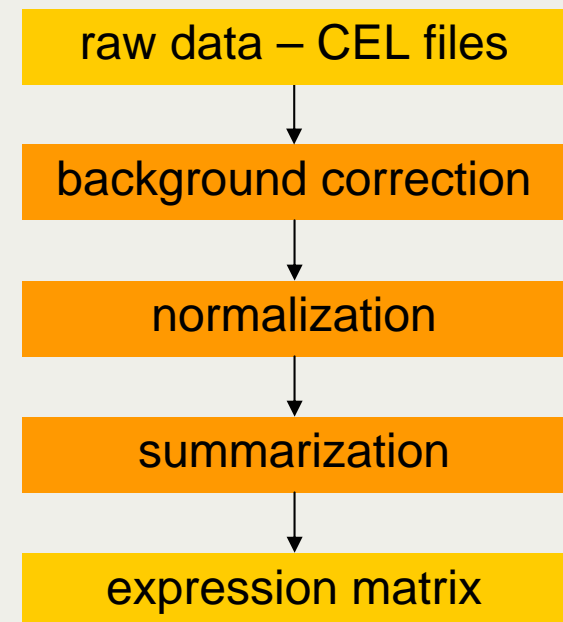- too random to be explicitely accounted for
- remain as "noise"

**Calibration**

**Error model**

4

# Preprocessing

- ## Background correction
  - remove noise of optical detection system

- ## Normalization
  - make measurements comparable from different array hybridizations

- ## Summarization
  - transcripts are represented in multiple probes

raw data – CEL files

↓

background correction

↓

normalization

↓

summarization

↓

expression matrix
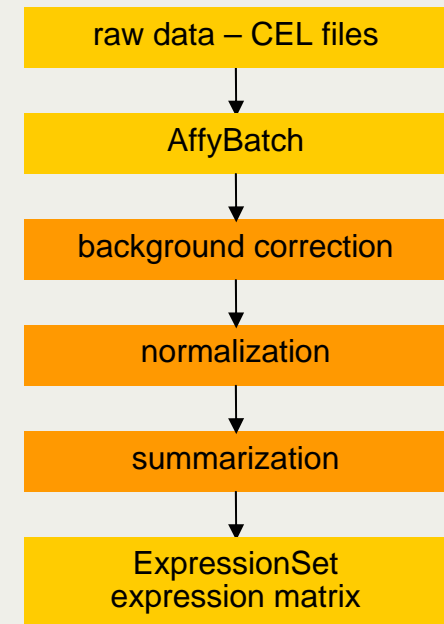
Markus Schmidberger, IBE
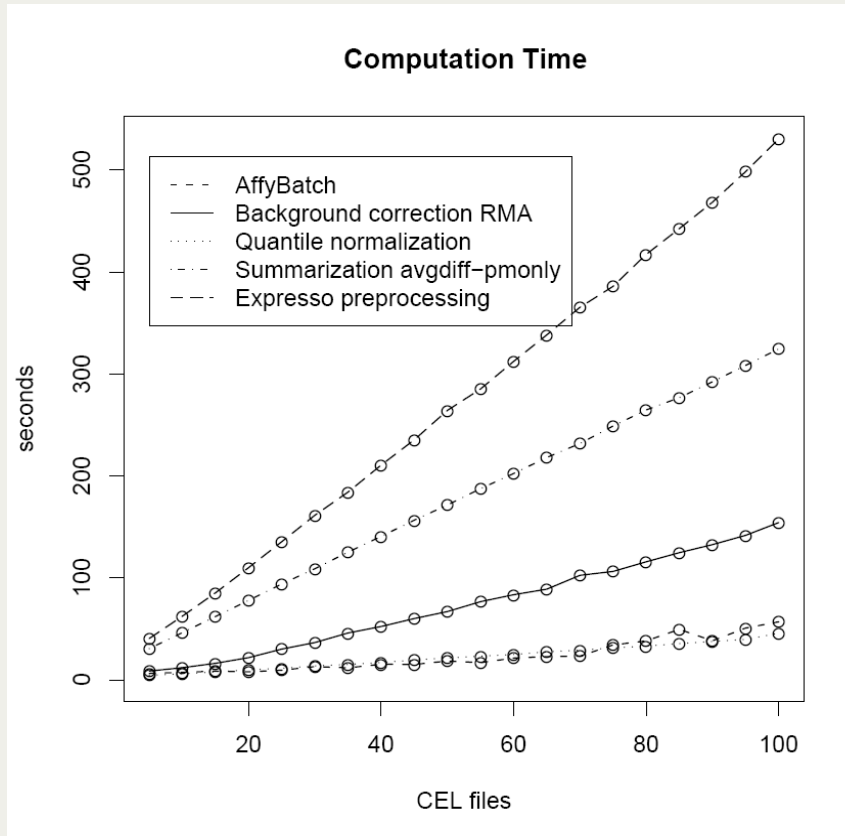HiCOMB 2008

# Existing software

- R and BioConductor mostly used in research (Open Source)

- Various different algorithms
  - With different advantages and disadvantages
  - No optimal solution (quality <-> effectiveness)

- Some approved and often used methods RMA, MAS 5.0, Quantile normalization, Cyclic loess, VSN, expresso, GCRMA

6

# Problems

- ## Data-structure of R
  - data are stored in class 'AffyBatch'
  - complex class with a lot of different slots
  - AffyBatch is memory intensiv

- ## Performance of algorithms
  - Inefficient program structure
  - Long computation time

raw data – CEL files

AffyBatch

background correction

normalization

summarization

ExpressionSet
expression matrix

7

# Problems

**Computation Time**

Legend:
- AffyBatch
- Background correction RMA
- Quantile normalization
- Summarization avgdiff–pmonly
- Expresso preprocessing

(y-axis: seconds; x-axis: CEL files)

**How many arrays can I RMA process?**
(Ben Bolstad)

http://bmbolstad.com/misc/Comp
uteRMAFAQ/size.html

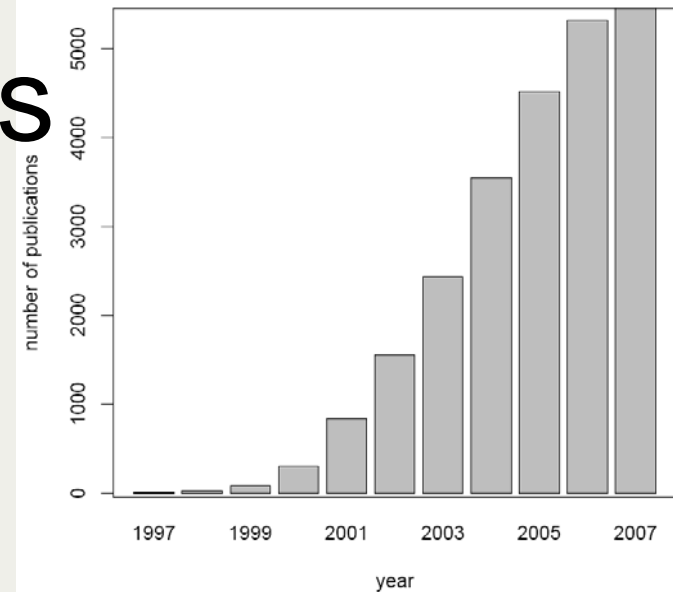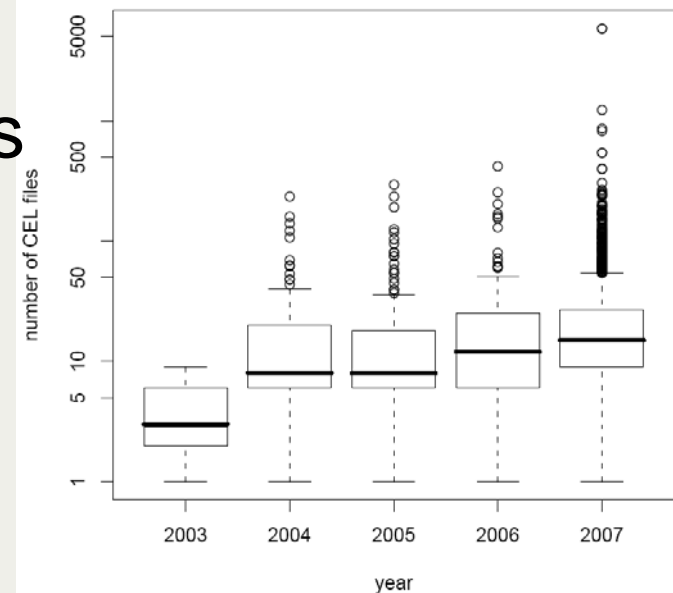| System | max. CEL files |
|---|---|
| 64-bit linux system with 4 GB main memory | 400 |
| 32-bit linux system with 4 GB main memory | 160 |
| 32-bit Microsoft Windows XP system with 1 GB main memory | 60 |

Chip: hgu133

# Challenges

- **Microarray experiments more and more popular**

- **Microarray chips become cheaper**
  - Experiments grow in size
  - EBI experiment: 6000 Arrays

- **Microarray chips grow in size**
  - More genes per chip



PubMed publications: 'microarray'



Size of experiments in ArrayExpress

# Possible Solutions

- Business applications
  - Expensive, not adaptable
- Faster and bigger computers
  - Expensive, limited
  - Main memory 256 GB: 60t €
- Better coding
  - C, DB
- Distribution to several computers / processors
  - Concurrent calculation of parts at different processors
  - Main memory 8 GB: 2000 € -> 60t € = 30 computers

10

# Parallelization

- ## Multiprocessors
  - the use of two or more central processing units (CPUs) within a single computer system
  - Today: Two-processors get a standard for workstations

- ## Multicomputers = Cluster
  - different parts of a program run simultaneously on two or more computers that are communicating with each other over a network
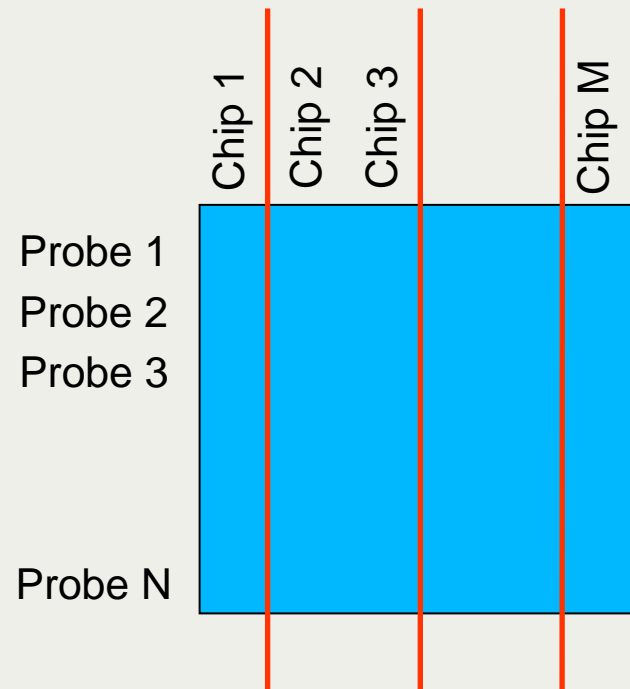  - Computer, network, software

# Software: MPI

- **Message Passing Interface**
- MPI is an API for parallel programming based on the message passing model
- MPI processes execute in parallel
- MPI is a standard for libraries
- Libraries exists for
  - FORTRAN, C, C++
  - R: Rmpi, Snow, papply

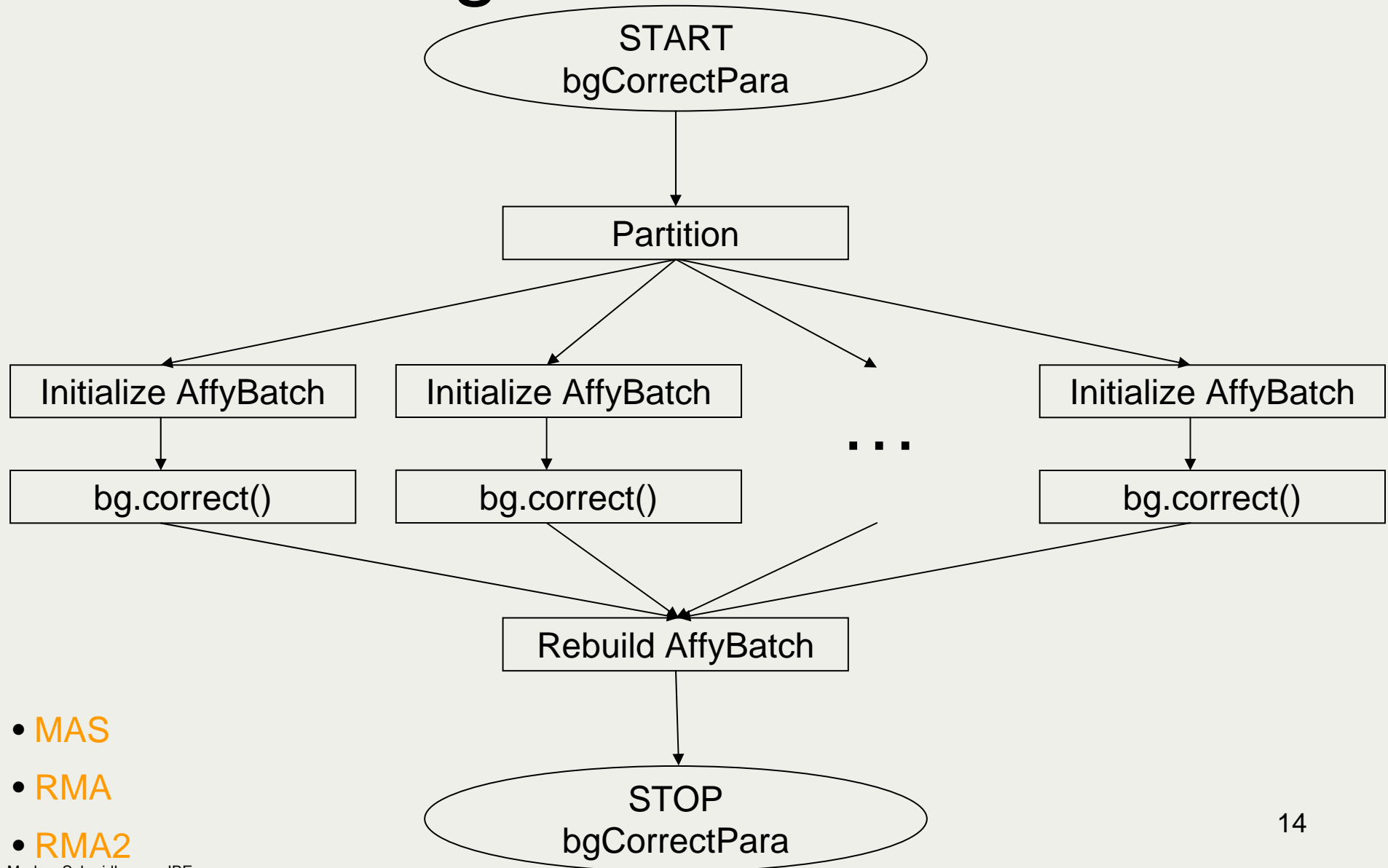- IBE Cluster: LAM/MPI 7.1.3



12

# Decomposition of AffyBatch

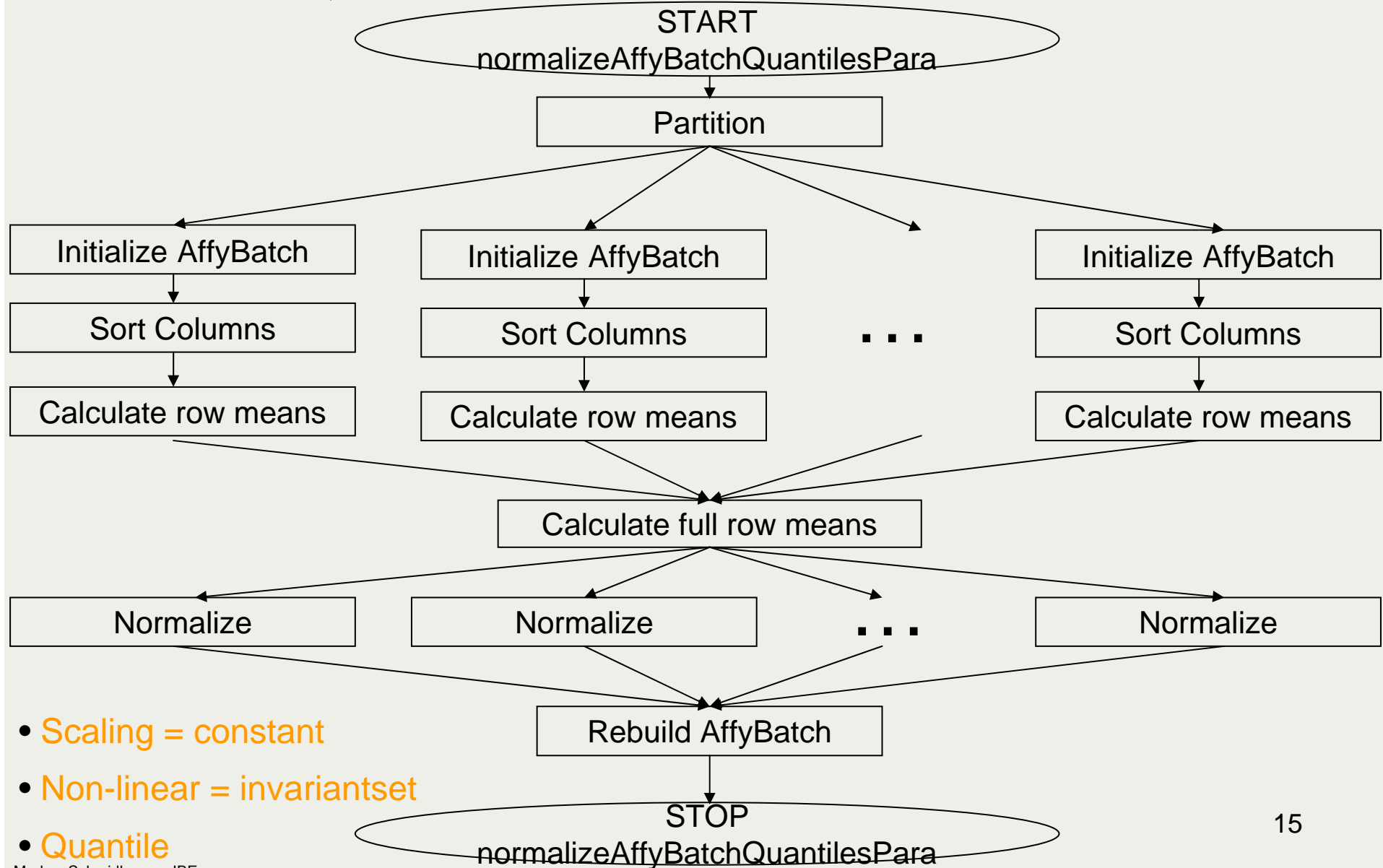- AffyBatch = intensities from multiple arrays



- Which decomposition is the best ?
  - Partition by chips
  - Partition by probes
  - Partition of CEL file name list
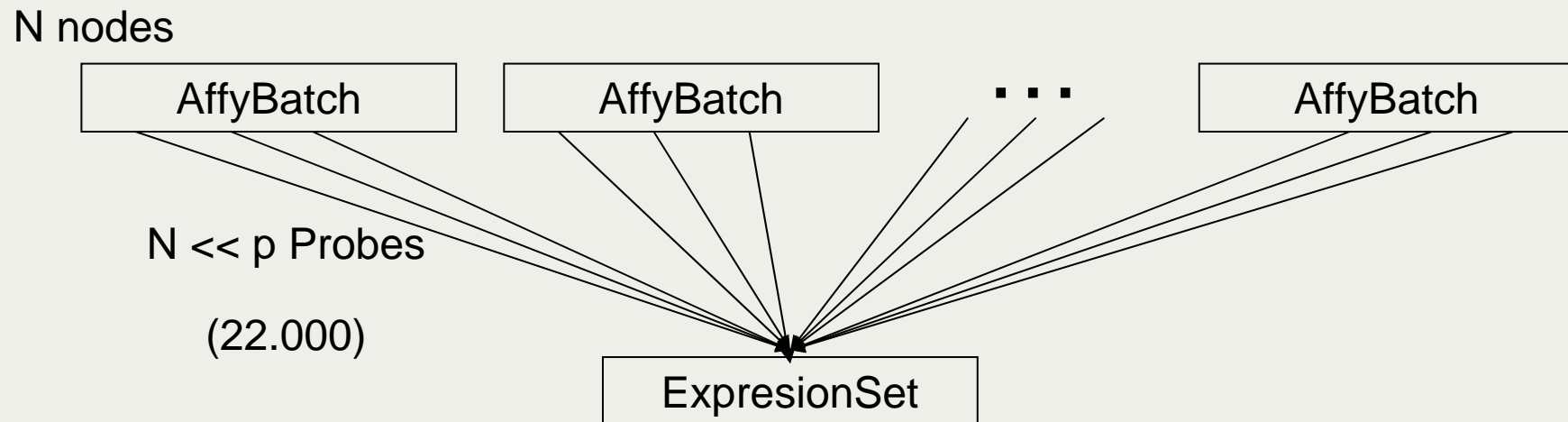
13

# Background Correction



- MAS
- RMA
- RMA2

# Quantile Normalization

```
         ┌─────────────────────────────────────┐
         │               START                 │
         │  normalizeAffyBatchQuantilesPara     │
         └─────────────────────────────────────┘
                          │
                  ┌───────────────┐
                  │   Partition   │
                  └───────────────┘
```

| Initialize AffyBatch | Initialize AffyBatch | | Initialize AffyBatch |
|---|---|---|---|
| Sort Columns | Sort Columns | . . . | Sort Columns |
| Calculate row means | Calculate row means | | Calculate row means |

```
          ┌──────────────────────────────┐
          │   Calculate full row means    │
          └──────────────────────────────┘
```

| Normalize | Normalize | . . . | Normalize |
|---|---|---|---|

```
          ┌──────────────────────────────┐
          │       Rebuild AffyBatch        │
          └──────────────────────────────┘
                          │
         ┌─────────────────────────────────────┐
         │               STOP                  │
         │  normalizeAffyBatchQuantilesPara     │
         └─────────────────────────────────────┘
```

- Scaling = constant

- Non-linear = invariantset

- Quantile

Markus Schmidberger, IBE
HiCOMB 2008

15

# Summarization

N nodes

| AffyBatch | AffyBatch | . . . | AffyBatch |
|-----------|-----------|-------|-----------|

N << p Probes

(22.000)

ExpresionSet

- avgdiff
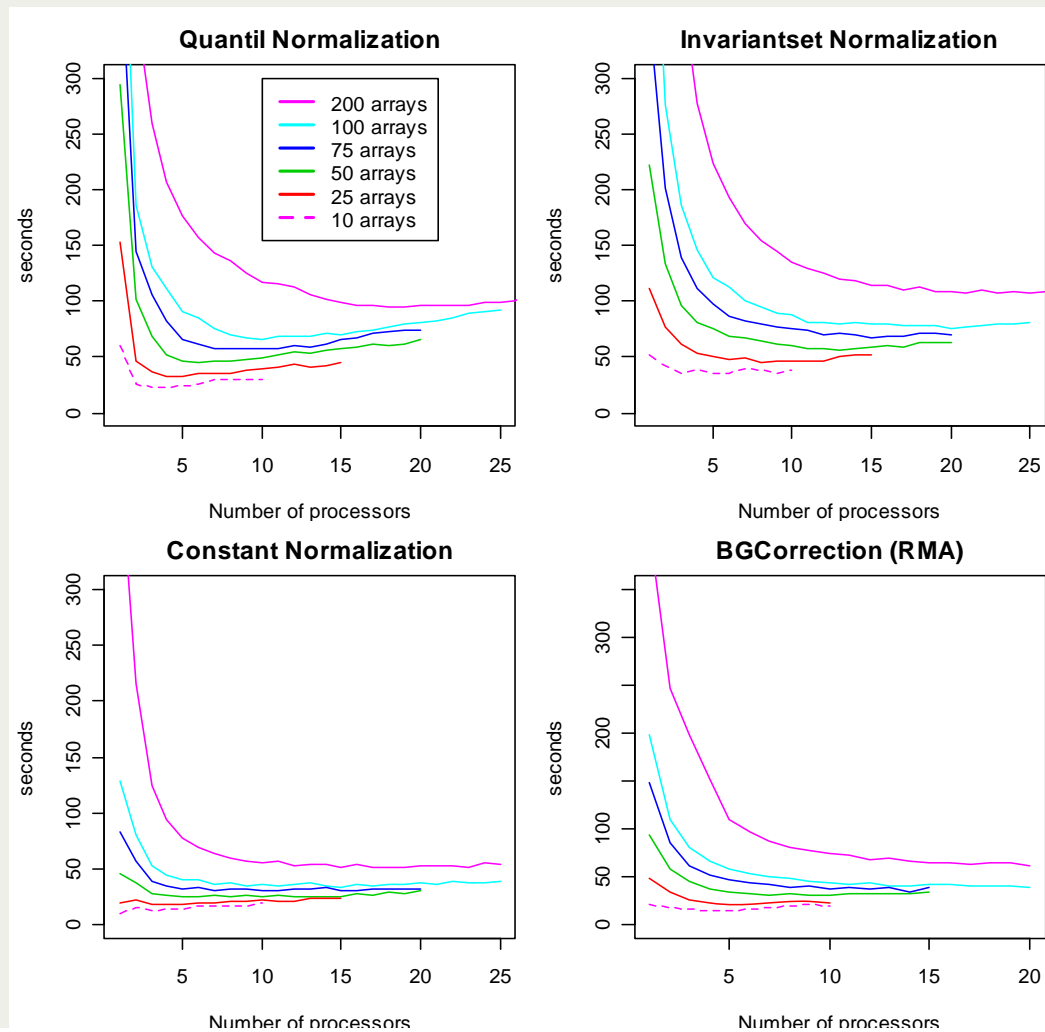- liwong
- mas
- medianpolish

# Communication Overhead

- **AffyBatch memory intensive**
  - A lot of data to transfer
  - Create AffyBatches at nodes
    - CEL Files available over samba device
  - Complete preprocessing method: preproPara()
    - Reducing the exchange of data
    - At no point a complete AffyBatch required

- **R functions and environments**
  - Define global functions
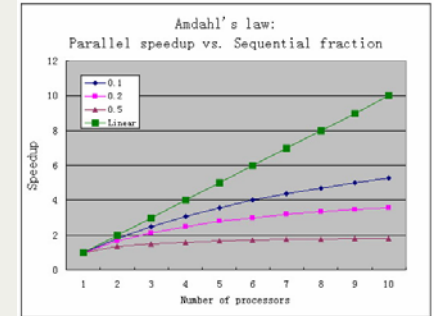
17

# Results

- Package **affyPara** with parallelized affy-functions
  - More CEL Files preprocessable
  - Speedup


- Parallelization methods produce in view of machine accuracy the <span style="color:orange">same results</span> as serialized methods.
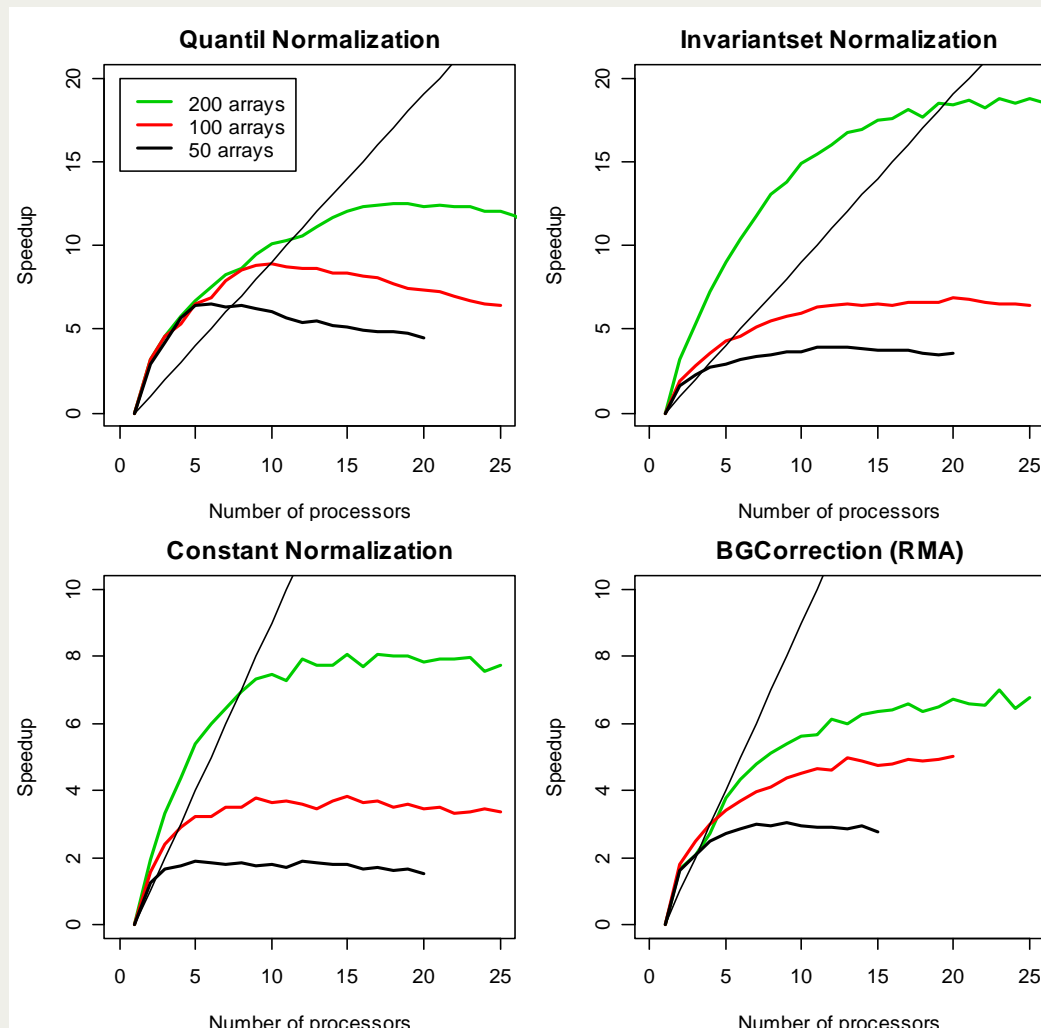  - All.equal(), machine's precision.

18

# Results - Computation Time

# Results – Speedup



Amdahl's law:
Parallel speedup vs. Sequential fraction

- **Speedup** of the methods up to factor 15

$Sp = T\_1 / T\_p$

$Sp \sim 1 / [ s + p/N ]$



Quantil Normalization

200 arrays
100 arrays
50 arrays

Speedup
Number of processors

Invariantset Normalization

Speedup
Number of processors

Constant Normalization

Speedup
Number of processors

BGCorrection (RMA)

Speedup
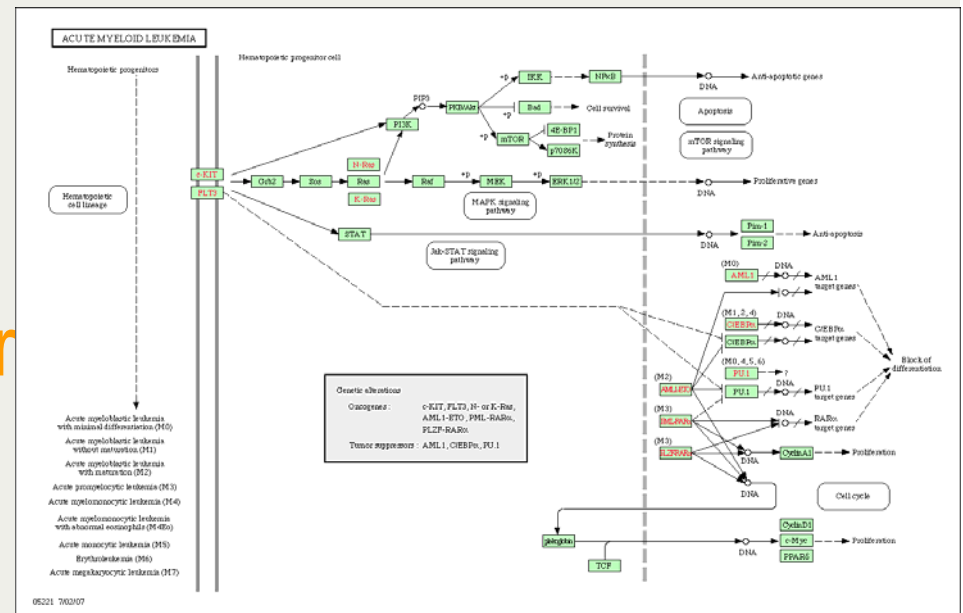Number of processors

20

# Results - Conclusion

- Partition of data and distribution to several nodes <span style="color:orange">solves the main memory problems</span>
  - IBE Cluster: ~ 16.000 microarrays
    - ( 32 workstations, 8 GB main memory, 2 dual core, Intel Xeon DP 5150, 1 Gbit Network, Linux 2.6.18 openSuse, LAM/MPI 7.1.3 )
  - Expansion of the cluster -> more data processible

- <span style="color:orange">affyPara</span> package available in the next BioConductor release in April 2008

# Large project in applied bioinformatics in preparation

- Collecting cancer data from public libraries
  - ArrayExpress, GEO, …
  - -> more than 5000 microarrays

- Preprocessing all together

- Analyzes all together
  - pathways ?



22

Markus Schmidberger, IBE
HiCOMB 2008

# Parallelized preprocessing algorithms for

# high-density oligonucleotide array data

**Thanks for your attention**

Dipl.-Tech. Math. Markus Schmidberger
schmidb@ibe.med.uni-muenchen.de

IBE
http://ibe.web.med.uni-muenchen.de

HiCOMB 2008
Seventh IEEE International Workshop on High
Performance Computational Biology
April 14, 2008; Miami, Florida, USA

23

# Literature R and Preprocessing

- Gentleman et. All; Bioinformatics and Computational Biology Solutions, *Using R and Bioconductor;* Springer, 2005 (Statistics for Biology and Health)

- Berrar et. All.; *A Practical Approach to Microarray Data Analysis;* Kluwer Academic Publishers, 2004

- Bolstad, Irizarry, Astrand, Speed; *A comparison of normalization methods for high density oligonucleotide array data based on variance and bias; Bioinformatics,* 2003, *19(2),* 185-193

- Irizarry, Wu, Jaffee; *Comparison of Affymetrix GeneChip expression measures; Bioinformatics, 2006, 22 no. 7,* 789-794

Markus Schmidberger, IBE
HiCOMB 2008

# Literature Parallel Computing

- Sloan, J. D.: *High Performance Linux Clusters with OSCAR, Rocks, OpenMosix, and MPI O'Reilly,* 2004
- Tierney, Luke: *Code Analysis and Parallelizing Vector Operations in R*. In: DSC 2007. Auckland, New Zealand, February 2007
- Rossini, Anthony: *Simple Parallel Statistical Computing in R*. In: UW Biostatistics Working Paper Series 193 (2003)
- Sevcikova, Hana: *Statistical Simulations on Parallel Computers*. In: Journal of Computational and Graphical Statistics 13 (2003), Nr. 4, 886-906.

Markus Schmidberger, IBE
HiCOMB 2008

# Package affyPara

| | | affy | affyPara |
|---|---|---|---|
| BGC | RMA | bg.correct.rma | bgCorrectPara |
| | MAS 5.0 | bg.correct.mas | bgCorrectPara |
| | RMA alt | bg.correct.rma2 | bgCorrectPara |
| Normalization | Scaling | normalize.AffyBatch.constant | normalizeAffyBatchConstantPara |
| | invariantset | normalize.AffyBatch.invariantset | normalizeAffyBatchInvariantsetPara |
| | Quantiles | normalize.AffyBatch.quantiles | normalizeAffyBatchQuantilesPara |
| Summarization | | computeExprSet | computeExprSetPara |
| complete | preprocessing | threestep, expresso | preproPara |