

# HPCS研究室 佐藤グループ(PA-team)の研究紹介

佐藤

筑波大学 計算科学研究センター

# 自己紹介



- <http://www.hpcs.cs.tsukuba.ac.jp/~msato/top.html.ja>
- 専門分野
  - 並列処理、並列システムの性能評価
  - 並列プログラミング言語、コンパイラ、OpenMP
  - グリッドコンピューティング、広域分散処理
  - 並列処理による低電力・省電力化
- 筑波大学 計算科学研究センター センター長
- 理研計算科学研究機構 プログラミング環境研究チーム・チームリーダー  
「京コンピュータ」の研究開発
- 公開ソフトウェア
  - Omni OpenMP/ XcalableMP コンパイラ <http://www.xcalablemp.org/>
  - OmniRPCシステム <http://www.omni.hpcc.jp/omniRPC/>
- 担当講義
  - プログラミング入門II(情報科学類)
  - プログラミング言語処理(情報科学類)
  - 分散システム(情報科学類)
  - プログラミング環境特論(システム情報工学研究科)

- 並列プログラミング言語XcalableMPをベースにした研究
  - コンパイラ
  - 大規模並列 – 性能プロファイル
  - GPU/メニーコア – HA-PACS
  - 耐故障 – 言語からのアプローチ
  
- マルチコア・メニーコアの並列プログラミングに関する研究
  - 細粒度並列処理技術 Cilkとか。
  - Work stealing scheduler ...
  
- 省電力化に関する研究
  - なんでも。

## ■ 現状と課題

- 並列プログラムの大半はMPI通信ライブラリによるプログラミング
  - 生産性が悪く、並列化のためのコストが高い。
- 並列プログラミングの教育のための簡便で標準的な言語がない(MPIでの教育にとどまっている)
- 研究室のPCクラスタから、センター、ペタコンまでに到るスケラブルかつポータブルな並列プログラミング言語が求められている

## ■ 目標

- 既存言語を指示文により拡張し、これからの大規模並列システム(分散メモリシステムと共有メモリノード)でのプログラミングを助け、生産性を向上させる並列プログラミング言語を設計・開発する。
- 標準化をすることを前提に、ユーザのわかりやすさを第一にどこでも使えるということを重視し、開発ならびに普及活動を進める。

### Current Problem ?!

```

int array[YMAX][XMAX];
main(int argc, char**argv){
  int i,j,res,temp_res, dx,llimit,ulimit,size,rank;

  MPI_Init(argc, argv);
  MPI_Comm_rank(MPI_COMM_WORLD, &rank);
  MPI_Comm_size(MPI_COMM_WORLD, &size);
  dx = YMAX/size;
  llimit = rank * dx;
  if(rank != (size - 1)) ulimit = llimit + dx;
  else ulimit = YMAX;

  temp_res = 0;
  for(i = llimit; i < ulimit; i++){
    for(j = 0; j < 10; j++){
      array[i][j] = func(i, j);
      temp_res += array[i][j];
    }
  }

  MPI_Allreduce(&temp_res, &res, 1, MPI_INT, MPI_SUM, MPI_COMM_WORLD);
  MPI_Finalize();
}
    
```

MPIしか、使えるものがない  
MPIの並列プログラムはむずかしい... いっぱい書き換えな  
いとけないし、時間がかかる、  
デバックもむずかしいし、...

### We need better solutions!!

```

#pragma xmp template T[10]
#pragma xmp distributed T[block]

int array[10][10];
#pragma xmp aligned array[i][*] to T[i]

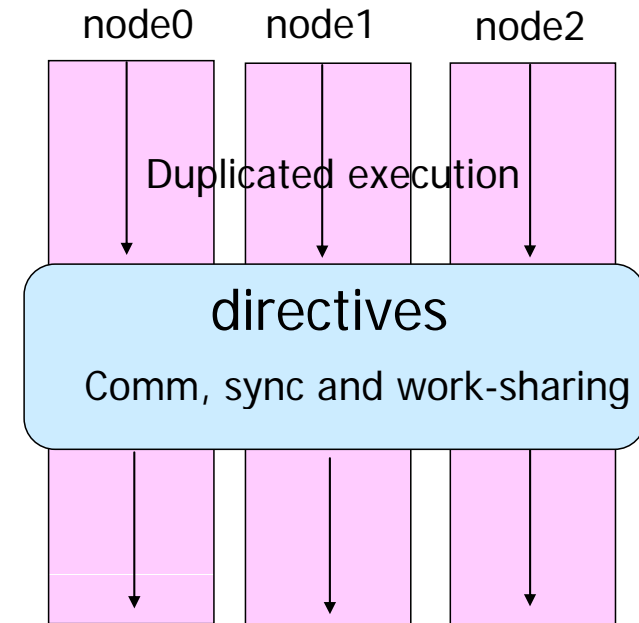
main(){
  int i, j, res;
  res = 0;
  #pragma xmp loop on T[i] reduction(+:
  for(i = 0; i < 10; i++){
    for(j = 0; j < 10; j++){
      array[i][j] = func(i, j);
      res += array[i][j];
    }
  }
}
    
```

data c  
add to  
increment  
work sharing and data  
synchronization

いまのプログラムに指示文を加える  
ただだから、簡単！性能チュー  
ニングも可能、...  
どこでも使えるから安心、並列プ  
ログラミングも習得にもお勧め！

## XcalableMP : directive-based language eXtension for Scalable and performance-tunable Parallel Programming

- **Directive-based language extensions** for familiar languages F90/C/C++
  - コードの書き換えや教育のコストを抑えること
- **“Scalable” for Distributed Memory Programming**
  - SPMDが基本的な実行モデル。
  - MPIのように、各ノードでスレッドが独立に実行を開始する。
  - 指示文(directive)がなければ、重複実行
  - タスク並列のためのMIMD実行も
- **“performance tunable” for explicit communication and synchronization.**
  - 指示文を実行するときに、Work-sharing や通信・同期がおきる。
  - すべての同期・通信操作は、指示文によって起きる。HPFと異なり、パフォーマンスのチューニングがわかりやすくなる。



# Code Example

```
int array[YMAX][XMAX];
```

```
#pragma xmp nodes p(4)  
#pragma xmp template t(YMAX)  
#pragma xmp distribute t(BLOCK) on p  
#pragma xmp align array[i][*] to t(i)
```

data distribution

```
main(){  
  int i, j, res;  
  res = 0;
```

add to the serial code : incremental parallelization

```
#pragma xmp loop on t[i] reduction(+:res)  
  for(i = 0; i < 10; i++)  
    for(j = 0; j < 10; j++){  
      array[i][j] = func(i, j);  
      res += array[i][j];  
    }  
}
```

work sharing and data synchronization

# XcalableMP コード例 (NPB CG, global view)



```
#pragma xmp nodes p[NPROCS]
#pragma xmp template t[N]
#pragma xmp distributed t[block] on p
...
#pragma xmp aligned [i] to t[i] :: x,z,p,q,r,w
#pragma xmp shadow [*] :: x,z,p,q,r,w
...
```

ノードの形状の定義

Templateの定義と  
データ分散を定義

データの分散は、  
templateにalign  
データの同期の  
ためのshadow  
を定義、この場  
合はfull shadow

Work sharing  
ループの分散

データの同期

```
/* code fragment from conj_grad in NPB CG */
sum = 0.0;
#pragma xmp loop on t[j] reduction(+:sum)
    for (j = 1; j <= lastcol-firstcol+1; j++) {
        sum = sum + r[j]*r[j];
    }
    rho = sum;
for (cgit = 1; cgit <= cgitmax; cgit++) {
#pragma xmp reflect p
#pragma xmp loop on t[j]
    for (j = 1; j <= lastrow-firstrow+1; j++) {
        sum = 0.0;
        for (k = rowstr[j]; k <= rowstr[j+1]-1; k++)
            sum = sum + a[k]*p[colidx[k]];
    }
    w[j] = sum;
}
#pragma xmp loop on t[j]
for (j = 1; j <= lastcol-firstcol+1; j++) {
    q[j] = w[j];
}
```

# HPCCベンチマークのプログラミングと性能



- HPC Challenge Benchmark Class2
  - 新しい並列プログラミング言語での記述性と性能を競うカテゴリ
    - Class1はシステム性能
  - 4つのベンチマーク
    - STREAM
    - Random Access
    - HPL
    - FFT
  - 今年は、Awardは、性能はIBM(X10 and UPC), 記述性はCary (Chapel) になった

SC09 HPCC Class2 でFinalist!



# Nicknamed the "K computer"



Kei (京) represents the numerical unit of 10 Peta ( $10^{16}$ ) in the Japanese language, representing the system's performance goal of 10 Petaflops. The Chinese character 京 can also be used to mean "a large gateway" so it could also be associated with the concept of a new gateway to computational science.

一、十、百、千、万、億、兆、京、垓、杼、穰、溝、澗、正、載、極、  
 $10^0$   $10^1$   $10^2$   $10^3$   $10^4$   $10^8$   $10^{12}$   $10^{16}$   $10^{20}$   $10^{24}$   $10^{28}$   $10^{32}$   $10^{36}$   $10^{40}$   $10^{44}$   $10^{48}$

恒河沙、阿僧祇、那由他、不可思議、無量大数  
 $10^{52}$   $10^{56}$   $10^{60}$   $10^{64}$   $10^{68}$

# K computer Delivery Began in Late September 2010

- The first eight racks of the K computer were delivered to Kobe from Fujitsu on September 28, 2010. More than 800 racks are required for a 10 Peta Flops Performance.
- A computer rack weighs about 1,300 kg in average. The rack contains 96 water-cooled Fujitsu SPARC64 VIIIfx CPU chips, each of which performs 128 GFlops, interconnected with the 3D Torus network that Fujitsu named Tofu.



Photo of First c

# Schedule of development



We are here.

		FY2006	FY2007	FY2008	FY2009	FY2010	FY2011	FY2012
<b>System</b>		Conceptual design		Detailed design		Prototype, evaluation	Production, installation, and adjustment	
								Tuning and improvement
<b>Applications</b>	Next-Generation Integrated Nanoscience Simulation	Development, production, and evaluation					Verification	
	Next-Generation Integrated Life Simulation	Development, production, and evaluation					Verification	
<b>Buildings</b>	Computer building	Design		Construction				
	Research building	Design		Construction				

K computer is on line.



AICS was founded in July 2010.

The computer building and research building are completed in May 2010

# RIKEN

## Advanced Institute for Computational Science (AICS)



- The institute have been established at the NGS in Kobe (started in October 2010)
- Missions
  - To run the K computer efficiently for users of wide research areas
  - Carry out the leading edge of computational science technologies and contribute for COE of computational science in Japan
  - Propose the future directions of HPC in Japan and conduct it.
- Topics
  - Promoting strong collaborations between computational and computer scientists, working with core-organizations of each fields together.
  - Fostering young scientists who exploit both computational and computer science
  - Research for new concepts for HPC in the future beyond the NGS (this is, exascale?)

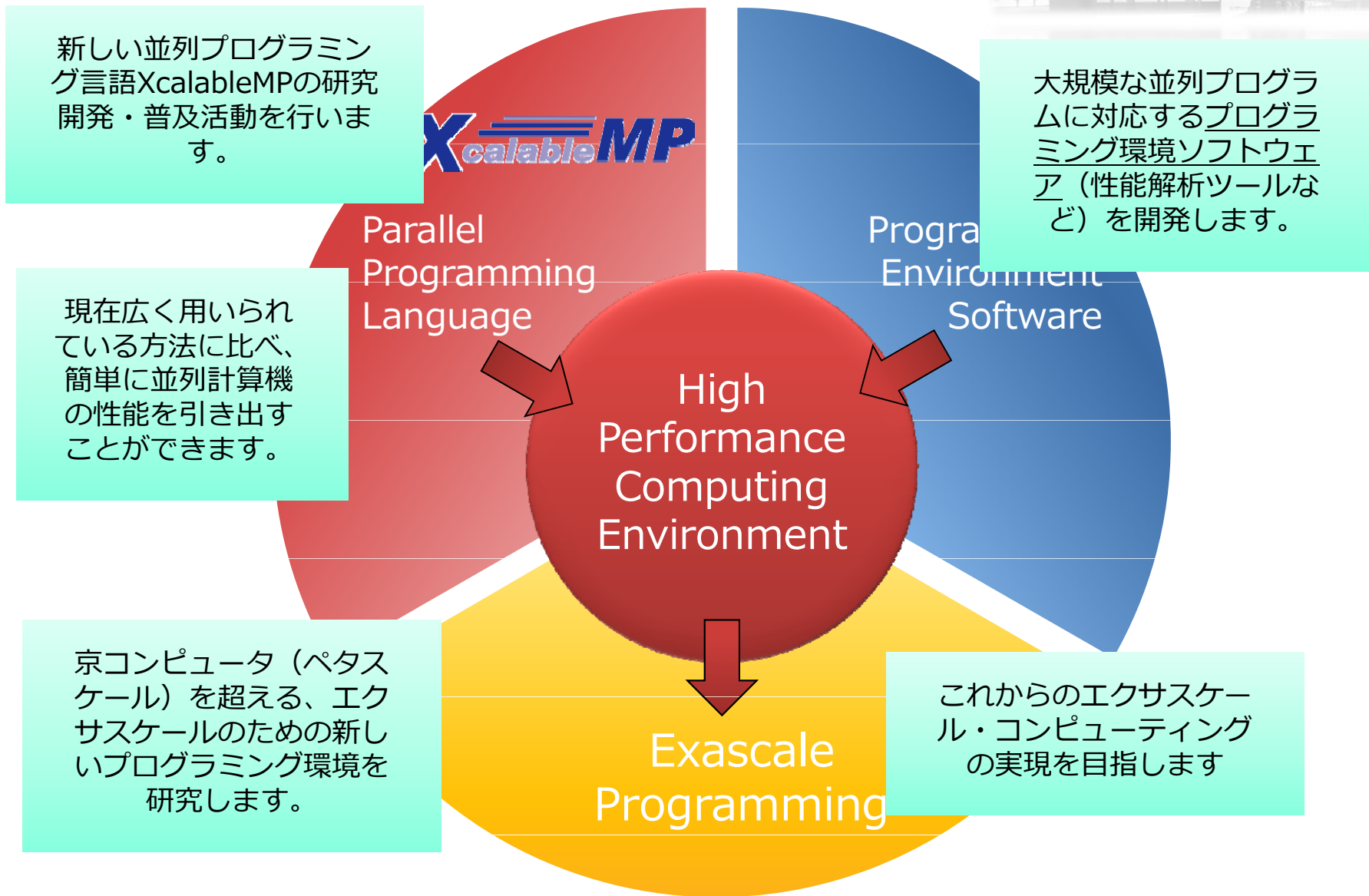
AICS

RIKEN Advanced Ins

<http://www.aics.riken.jp/>



# プログラミング環境研究チームの研究概要



- 並列プログラミング言語XcalableMPをベースにした研究
  - コンパイラ
  - 大規模並列 – 性能プロファイル
  - GPU/メニーコア – HA-PACS
  - 耐故障 – 言語からのアプローチ
  - フランスとの共同研究(フランスに行きたいひと！)
- マルチコア・メニーコアの並列プログラミングに関する研究
  - 細粒度並列処理技術 Cilkとか。
  - Work stealing scheduler ...
- 省電力化に関する研究
  - 応、相談