

## Stéphane Bihan, CAPS





#### Introduction

- Main stream applications will rely on new multicore / manycore architectures
  - It is about performance not parallelism
- Various heterogeneous hardware
  - General purpose cores
  - Application specific cores GPUs (HWAs)
- HPC and embedded applications are increasingly sharing characteristics



# An Overview of Hybrid Parallel Computing



### Manycore Architectures

- General purpose cores
  - Share a main memory
  - Core ISA provides fast SIMD instructions
- Streaming engines / DSP / FPGA
  - Application specific architectures ("narrow band")
  - Vector/SIMD
  - Can be extremely fast
- Hundreds of GigaOps
  - But not easy to take advantage of
  - One platform type cannot satisfy everyone
- Operation/Watt is the efficiency scale







#### Multicore/Manycore Workload

#### Multiple applications sharing the hardware

- Multimedia, game, encryption, security, health, ...
- Unfriendly environment with many competitions
  - Global resource allocation, no warranty on availability
  - Must be taken into account when programming/compiling
- Applications cannot always be recompiled
  - Most applications are distributed as binaries
- A binary will have to run on many platforms
  - Forward scalability or "write once, run faster on new hardware"
  - Loosing performance is not an option





- Amdahl's law is forever, all levels of parallelism need to be exploited
- Programming various hardware components of a node cannot be done separately



The Past of Parallel Computing,



the Future of Manycores?

#### The Past

- Scientific computing focused
- Microprocessor or vector based, homogeneous architectures
- Trained programmers willing to pay effort for performance
- Fixed execution environments
- The Future
  - New applications (multimedia, medical, ...)
  - Thousands of heterogeneous systems configurations
  - Unfriendly execution environments



#### Programming Multicores/

#### Manycores

- Physical architecture oriented
  - Shared memory architectures
    - OpenMP, CILK, TBB, automatic parallelization, vectorization...
  - Distributed memory architectures
    - Message passing, PGAS (Partition Global Address Space), ...
  - Hardware accelerators, GPU
    - CUDA, OpenCL, Brook+, HMPP, ...
- Different styles
  - Libraries
    - MPI, pthread, TBB, SSE intrinsic functions, ...
  - Directives
    - OpenMP, HMPP, ...
  - Language constructs
    - UPC, Cilk, Co-array Fortran, UPC, Fortress, Titanium, ...



Parallel Hybrid Computing, Euro GPU 2009, Lyon, Sept. 1 2009

#### Multi (languages) programming

- Happens when programmers need to deal with multiple programming languages
  - E.g. Fortran and CUDA, Java and OpenCL, ...
- Multiprogramming impacts on
  - Programmer's expertise
  - Program maintenance and correctness
  - Long term technology availability
- Performance programming versus domain specific programming
  - Libraries, parallel components to be provided to divide the issues



#### Manycore =



#### Multiple µ-Architectures

- Each µ–architecture requires different code generation/ optimization strategies
  - Not one compiler in many cases
- High performance variance between implementations
  - ILP, GPCore/TLP, HWA
- Dramatic effect of tuning
  - Bad decisions have a strong effect on performance
  - Efficiency is very input parameter dependent
  - Data transfers for HWA add a lot of overheads

### How to organize the compilation flow?









#### Can the Hardware be Hidden?

- Programming style is usually hardware independent but
  - Programmers need to take into account available hardware resources
- Quantitative decisions as important as parallel programming
  - Performance is about quantity
  - Tuning is specific to a configuration
- Runtime adaptation is a key feature
  - Algorithm, implementation choice
  - Programming/computing decision



- Available hardware resources are changing over the execution time
  - Not all resources are time-shared, e.g. a HWA may not be available
  - Data affinity must be respected

How to ensure that conflicts in resource usage will not lead to global performance degradation?





- Decision is a complex issue
  - How to produce the decision?





#### **Research Directions**

- New Languages
  - X10, Fortress, Chapel, PGAS languages, OpenCL, MS Axum, ...
- Libraries
  - Atlas, MKL, Global Array, Spiral, Telescoping languages, TBB, ...
- Compilers
  - Classical compiler flow needs to be revisited
  - Acknowledge lack of static performance model
  - Adaptative code generation
- OS
  - Virtualization/hypervisors
- Architectures
  - Integration on the chip of the accelerators
    - AMD Fusion, ...
  - Alleviate data transfers costs
    - PCI Gen 3x, ...



Key for the short/mid term

# HMPP Approach



#### HMPP Objectives

- Efficiently orchestrate CPU/GPU computations in legacy code
  - With OpenMP-like directives
- Automatically produce tunable manycore applications
  - C and Fortran to CUDA data parallel code generator
  - Make use of available compilers to produce binary
- Ease application deployment

#### HMPP....

#### a high level abstraction for manycore programming





#### HMPPI.5 Simple Example

```
#pragma hmpp label codelet, target=CUDA:BROOK, args[v1].io=out
#pragma hmpp label2 codelet, target=SSE, args[v1].io=out, cond="n<800"
void MyCodelet(int n, float v1[n], float v2[n], float v3[n])
{ int i;
   for (i = 0 ; i < n ; i++) {
     v1[i] = v2[i] + v3[i];
   }
}
Parallel Hybrid Computing, Euro GPU 2009, Lyon, Sept. 1 2009</pre>
```



#### Group of Codelets (HMPP 2.0)

- Declare group of codelets to optimize data transfers
- Codelets can share variables
  - Keep data in GPUs between two codelets
  - Avoid useless data transfers
  - Map arguments of different functions in same GPU memory location (equivalence Fortran declaration)

#### Flexibility and Performance





#### Optimi∠ing Communications

#### Exploit two properties

- Communication / computation overlap
- Temporal locality of parameters
- Various techniques
  - Advancedload and Delegatedstore
  - Constant parameter
  - Resident data
  - Actual argument mapping





### Hiding Data Transfers

- Pipeline GPU kernel execution with data transfers
  - Split single function call in two codelets (C1, C2)







#### Advancedload Directive

#### Avoid reloading constant data



t2 is not reloaded each loop iteration





#### Actual Argument Mapping

- Allocate arguments of various codelets to the same memory space
  - Allow to exploit reuses of argument to reduce communications
  - Close to equivalence in Fortran



#### HMPP Tuning

```
!$HMPP sgemm3 codelet, target=CUDA, args[vout].io=inout
SUBROUTINE sgemm(m,n,k2,alpha,vin1,vin2,beta,vout)
INTEGER, INTENT(IN) :: m,n,k2
REAL,
                     :: alpha, beta
       INTENT(IN)
REAL,
       INTENT(IN) :: vin1(n,n), vin2(n,n)
REAL, INTENT(INOUT) :: vout(n,n)
REAL
       :: prod
                                                    X>8 GPU compiler fails
INTEGER :: i,j,k
!$HMPPCG unroll(X), jam(2), noremainder
!$HMPPCG parallel
                                                    X=8 200 Gigaflops
DO j=1,n
    !$HMPPCG unroll(X), splitted, noremainder
                                                    X=4 100 Gigaflops
    !$HMPPCG parallel
   DO i=1,n
       prod = 0.0
       DO k=1,n
         prod = prod + vin1(i,k) * vin2(k,j)
        ENDDO
        vout(i,j) = alpha * prod + beta * vout(i,j) ;
     END DO
END DO
END SUBROUTINE sgemm
```



## -

### Conclusion

- Multicore ubiquity is going to have a large impact on software industry
  - New applications but many new issues
- Will one parallel model fit all?
  - Surely not but multi languages programming should be avoided
  - Directive based programming is a safe approach
  - Ideally OpenMP will be extended to HWA
- Toward Adaptative Parallel Programming
  - Compiler alone cannot solve it
  - Compiler must interact with the runtime environment
  - Programming must help expressing global strategies / patterns
  - Compiler as provider of basic implementations
  - Offline-Online compilation has to be revisited

