

并行程序设计原理

北京大学信息学院



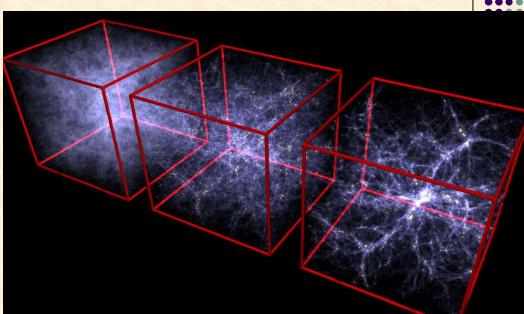
内容

- 并行体系结构
- 低层并行编程介绍 (CUDA、MPI、Pthread、SCIF、GlobalArrays、OpenMP Offload)
- 典型并行算法
- 高级并行编程 (OpenMP、Parray、MapReduce)
- 并行理论 (Petri网、线性逻辑、进程代数)

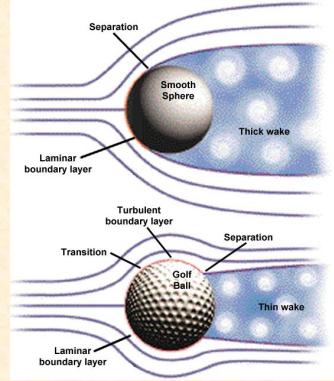
Getting Started

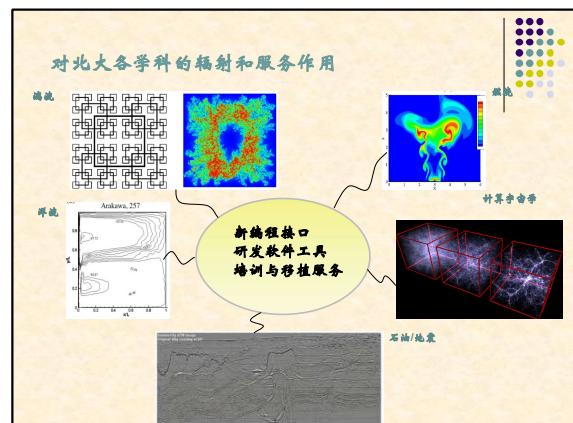
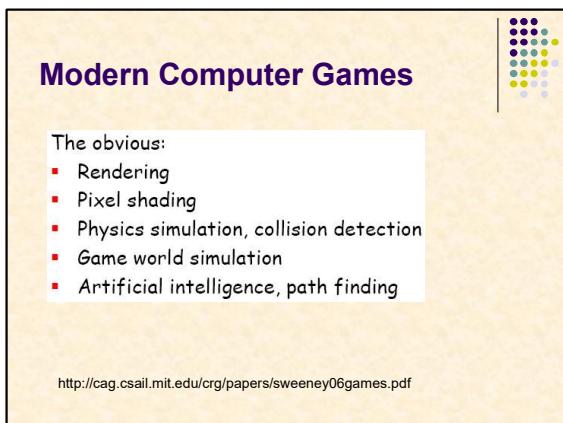
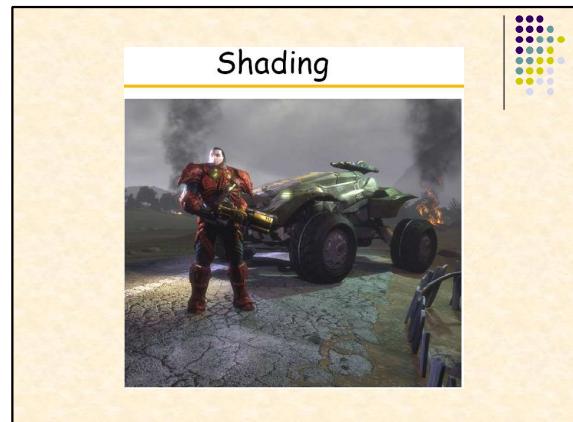
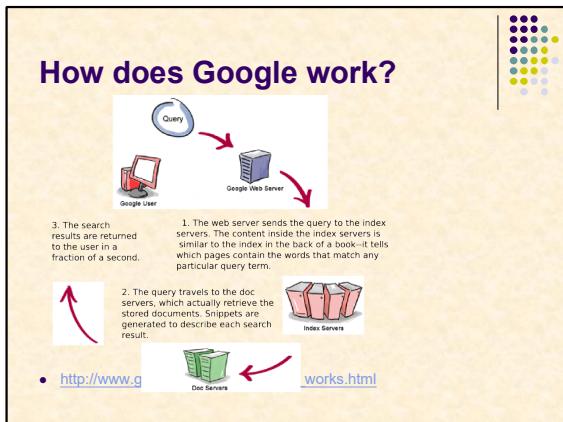
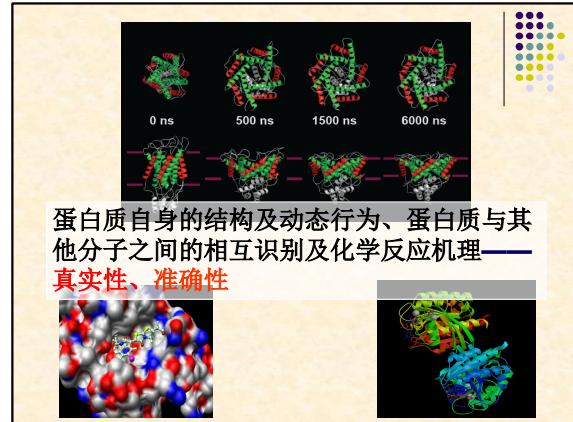
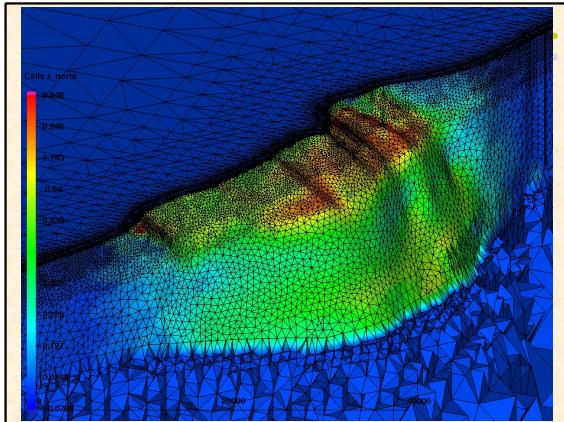
- http://www.nvidia.com/object/cuda_get.html

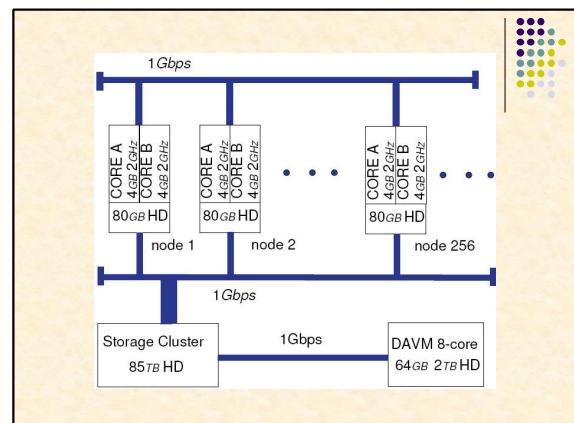
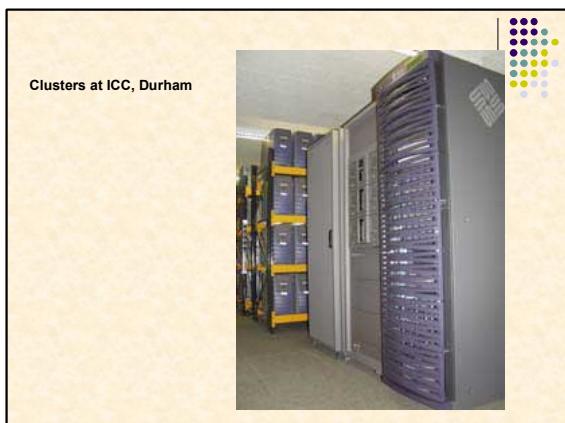
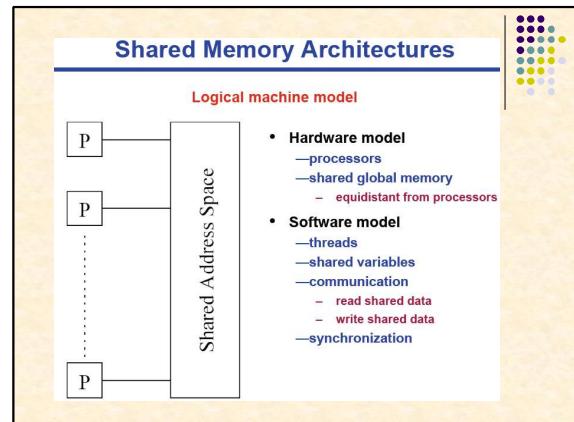
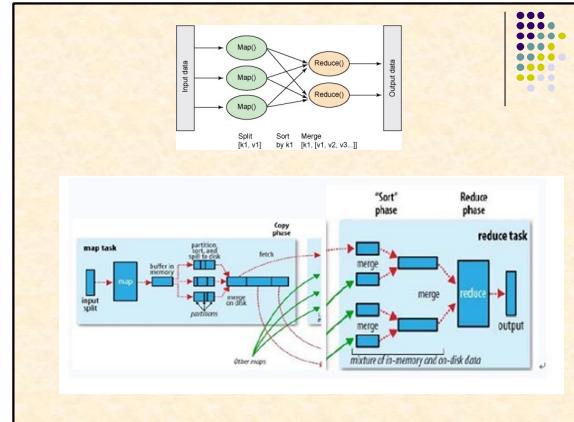
并行计算的应用

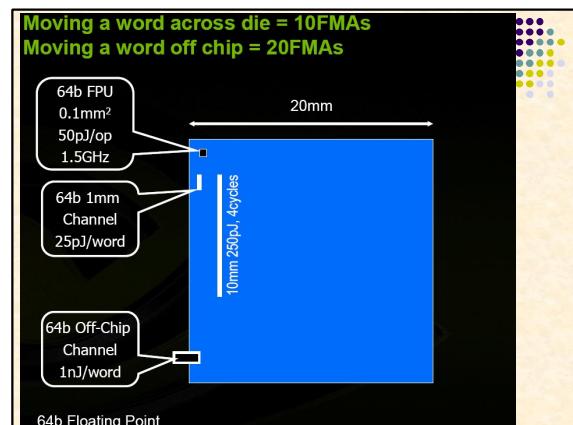
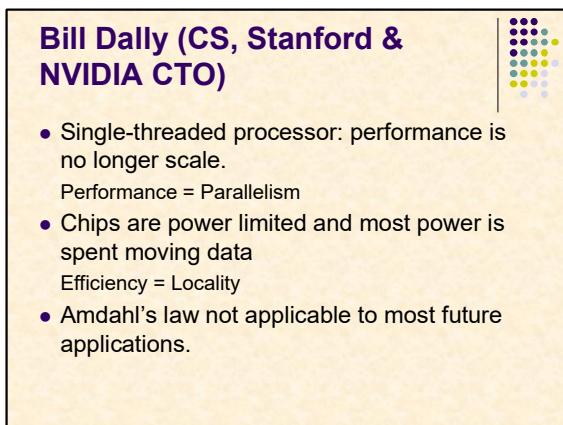
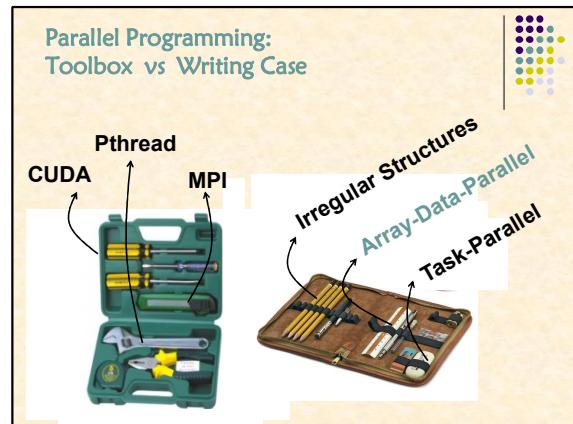
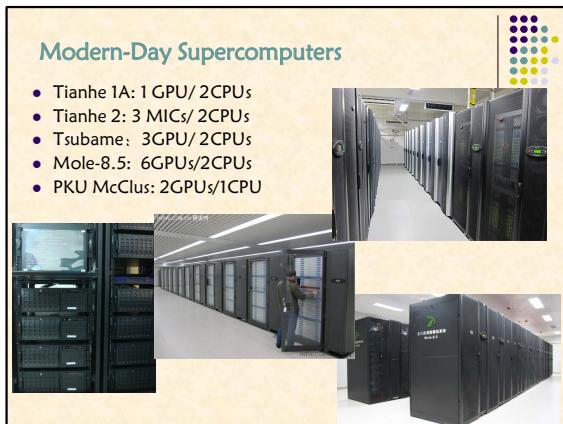
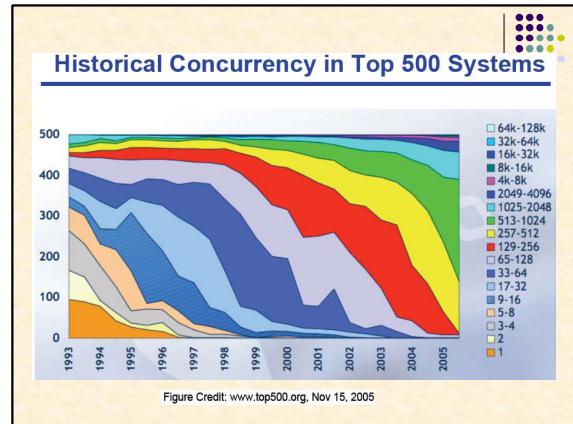
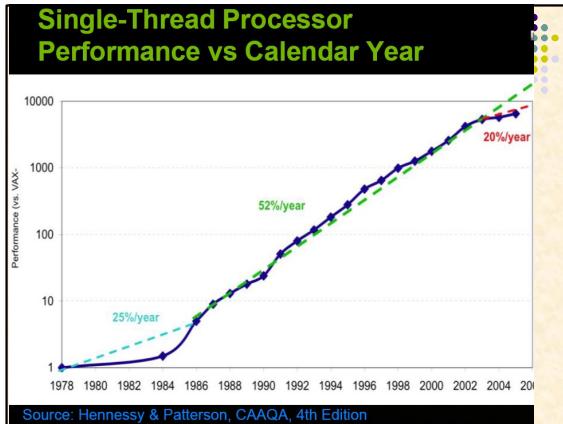



cosmological numerical simulations Springel et al. (2001)



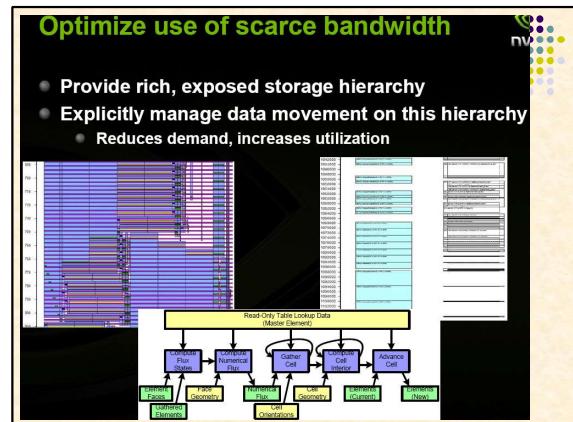
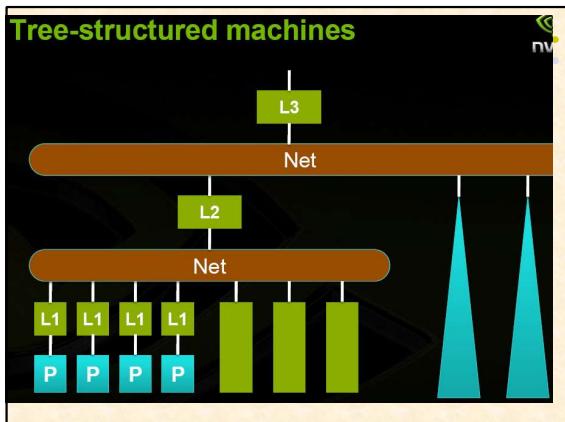
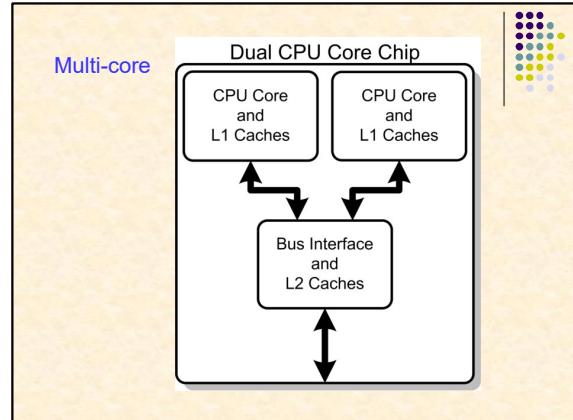







Exploiting parallelism and locality requires:

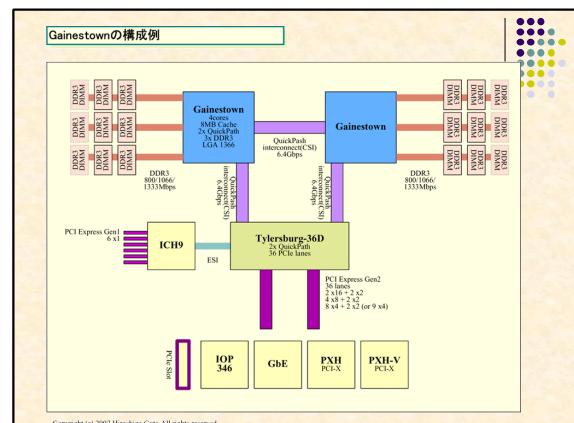
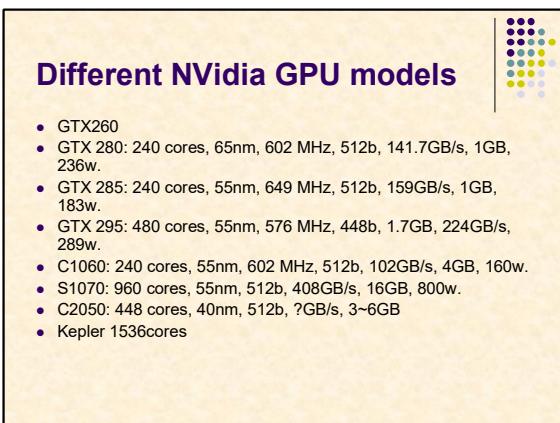
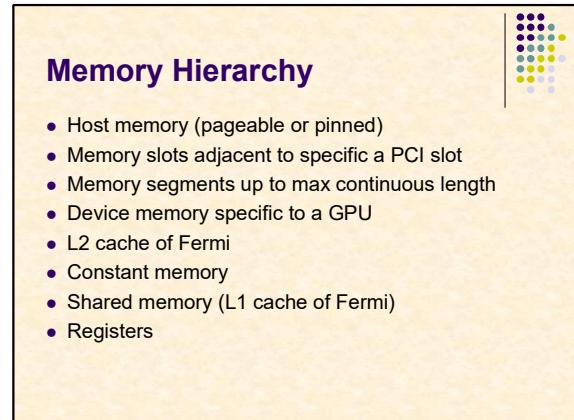
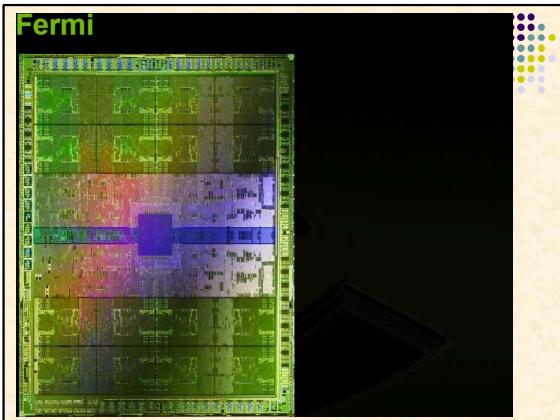
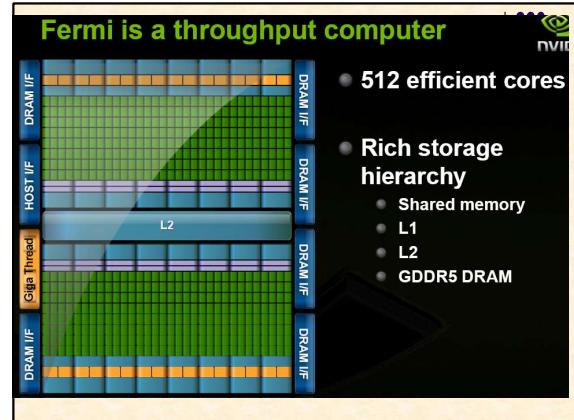
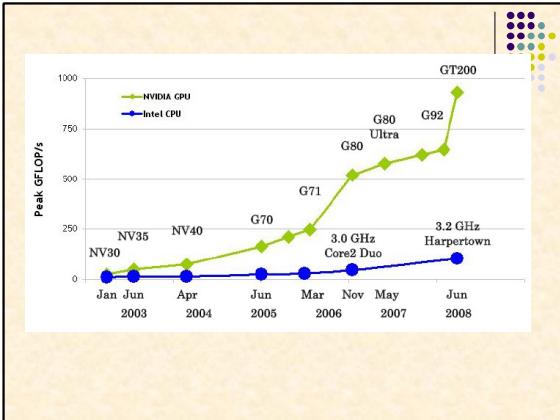
- Many efficient processors**
(To exploit parallelism)
- An exposed storage hierarchy**
(To exploit locality)
- A programming system that abstracts this**

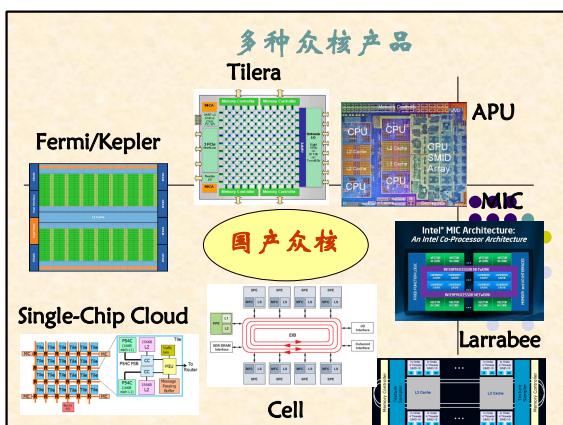
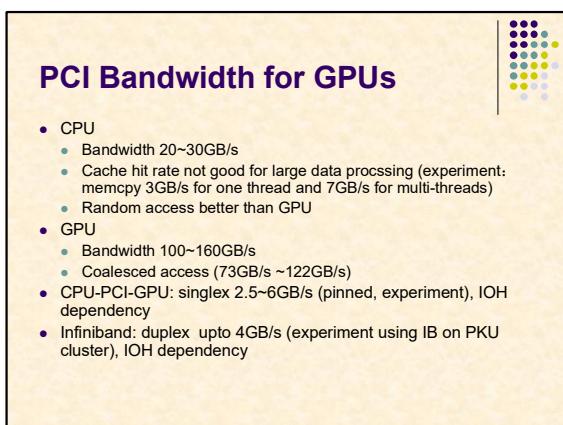
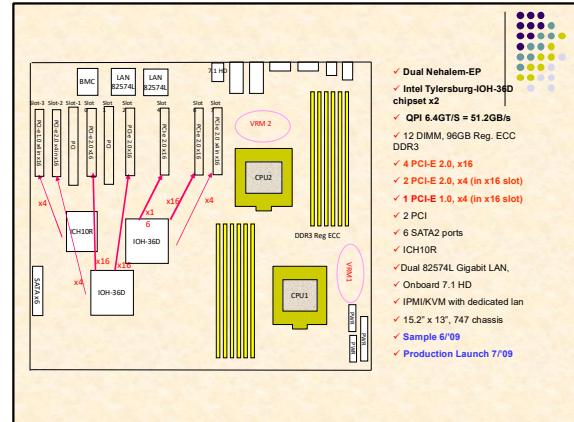
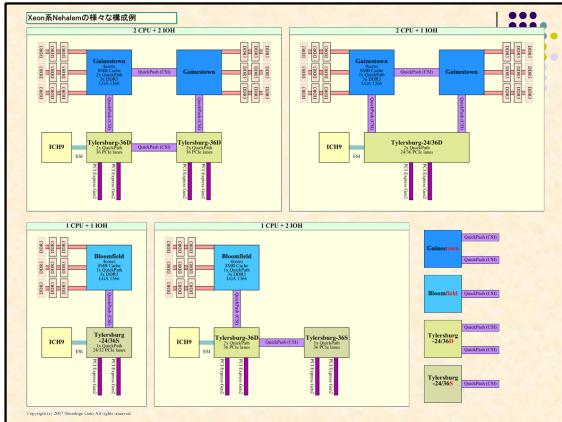


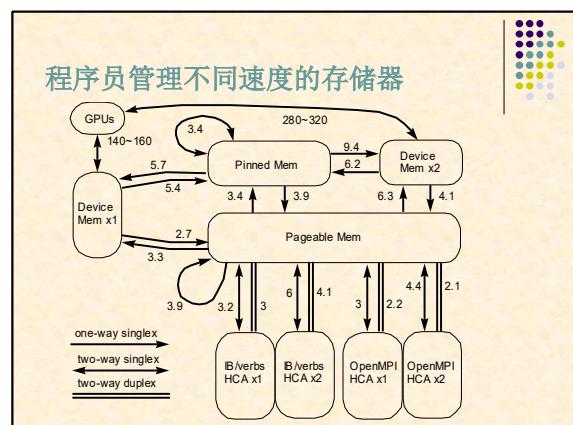
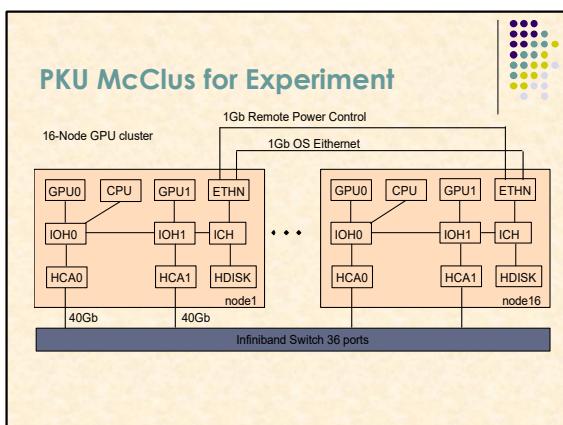
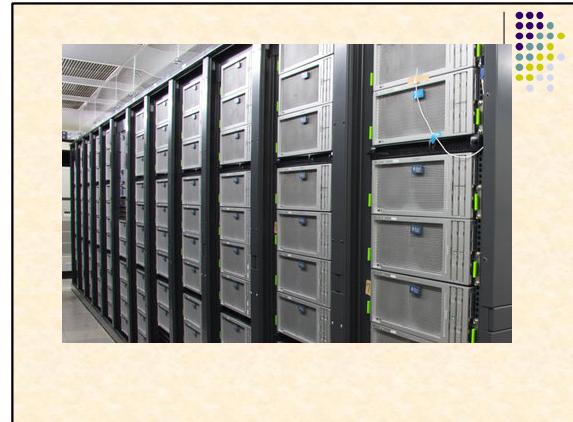
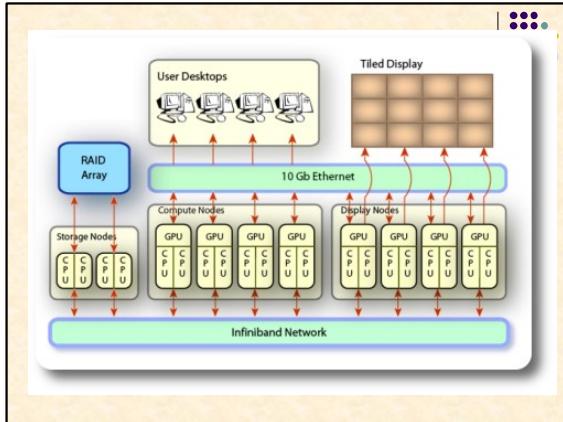
Avoid Denial Architecture

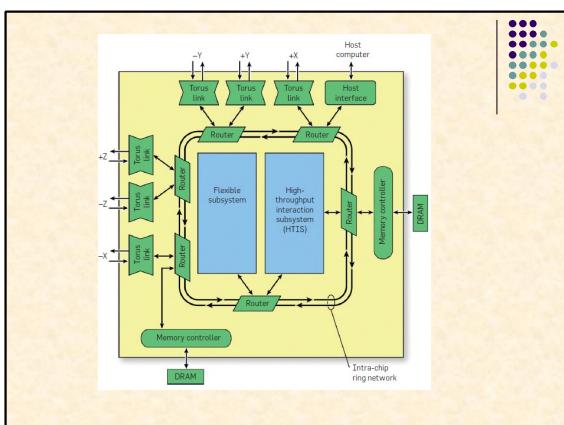
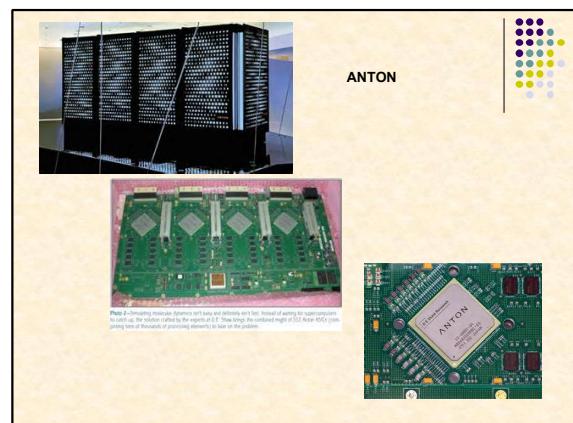
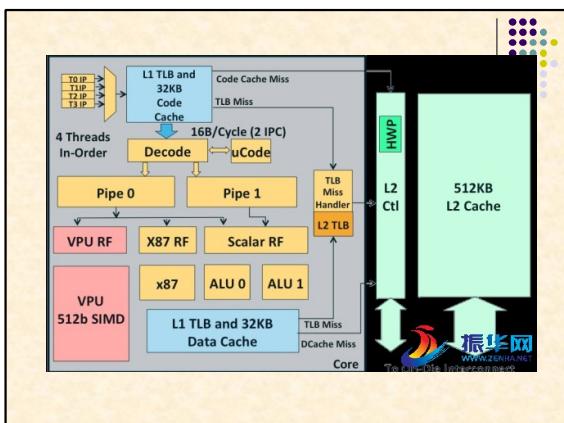
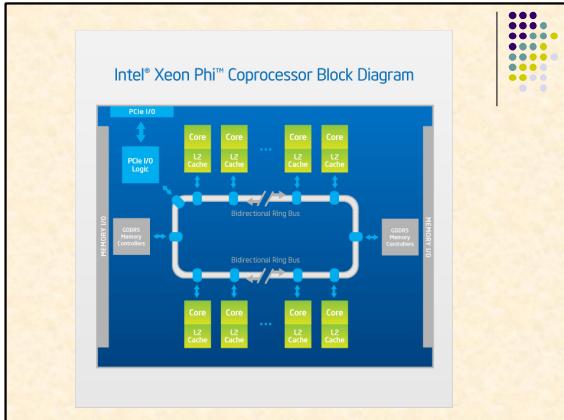
- Single thread processors are in denial about parallelism and locality
- They provide two illusions:
 - Serial execution - Denies parallelism
 - Tries to exploit parallelism with ILP - inefficient & limited scalability
 - Flat memory - Denies locality
 - Tries to provide illusion with caches – very inefficient when working set doesn't fit in the cache
- These illusions inhibit performance and efficiency







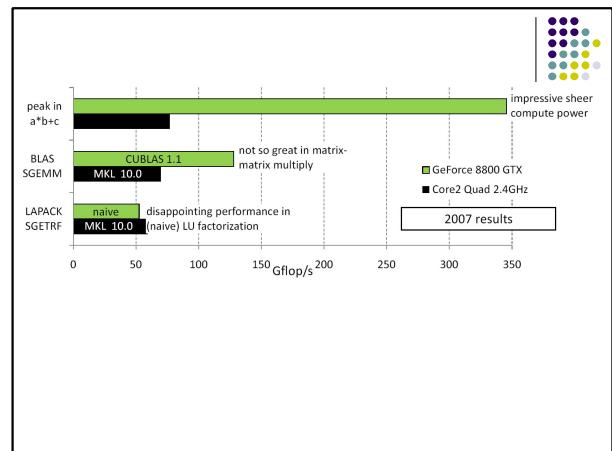




SGEMM

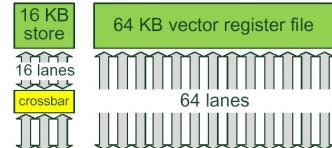
- Vasily Volkov, James Demmel:
Benchmarking GPUs to tune dense linear algebra. SC 2008

Some following slides are from the presentation of the above paper at SC'08.



GPU (NVIDIA GeForce)	8600 GTS	9800 GTX	GTX 280
Processor cores	4	16	30
Compute capability ($a+b+c$)	93 Gflop/s	429 Gflop/s	624 Gflop/s
Memory controllers	2	4	8
Memory bandwidth	32 GB/s	70 GB/s	141 GB/s

GPU Memory Hierarchy



- Register file is the fastest and the largest on-chip memory
 - Keep as much data as possible in registers
 - However, register file is constrained to vector operations
 - Can live with it — vectorized codes are common in HPC
- Shared memory permits indexed and shared access
 - However, it is 2–4x smaller and has 4x lower bandwidth than registers
 - Only 1 operand in shared memory is allowed versus 4 register operands
 - Moreover, some instructions run slower if using shared memory
 - Use shared memory as a communication device only
 - Avoid communication to improve performance

Peak Throughput in Multiply-and-Add

- How much parallelism is enough to get the peak?
- Run 1 thread per processor core
 - Purpose: smallest amount that can control all computing resources
- Assume sufficient instruction-level parallelism in the program
 - Purpose: hide pipeline latency
- Choose the shortest vector length that yields the peak
 - Purpose: satisfy inherent data-parallelism constraints
- Result: 98% of arithmetic peak at VL = 64
 - Therefore, VL=64 is recommended for all compute-bound codes
- However, we never could surpass 66% of peak is using an operand in shared memory
 - We believe this is an inherent bottleneck in the architecture
 - We use this number in the throughput bounds below

Matrix-Matrix Multiply: $C = C + A^*B$

- GPU requires using block algorithms in matrix-matrix multiply:
 - Peak rates on one of the latest GPUs are 624 Gflop/s and 141 GB/s
 - This corresponds to 0.23 bytes per flop
 - But naïve matrix-matrix multiply requires 4 bytes per flop
 - Thus, it is bandwidth-bound unless data is reused 18 times
 - Using $M \times N$ blocks in C yields $2/(1/M+1/N)$ average reuse
- Use vector algorithms to efficiently use vector registers
 - Such as used on IBM 3090 Vector Facility and Cray X1:
 - Keep A 's and C 's blocks in registers
 - Keep B 's block in a shared storage
 - No other sharing is needed if C 's height = VL. We know VL=64 is best
- Choose large enough width of C 's block
 - 16 is enough as $2/(1/64+1/16) = 26$ -way reuse
- Choose a convenient thickness for A 's and B 's blocks

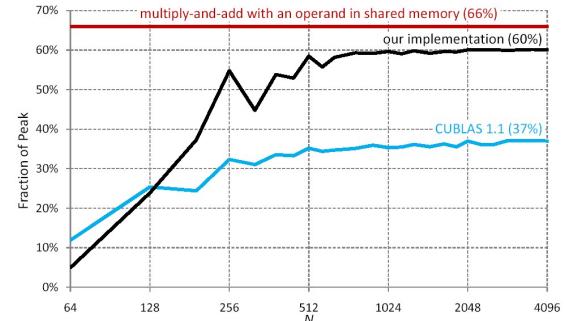
```

__global__ void sgemmNN(const float *A, int lda, const float *B, int ldb, float *C, int ldc, int k, float alpha, float beta)
{
    A += blockIdx.x * 64 + threadIdx.x * 16;
    B += threadIdx.x * 16 + threadIdx.y * ldb;
    C += blockIdx.x * 64 + threadIdx.x * (threadIdx.y + blockIdx.y * ldc) * 16; } Compute pointers to the data
__shared__ float bs[16][17];
float c[16] = {0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0}; } Declare the on-chip storage
const float *Blast = B + k;
do
{
#pragma unroll
    for( int i = 0; i < 16; i+= 4 )
        bs[threadIdx.x][threadIdx.y+i] = B[i*ldb]; } Read next B's block
    B += 16;
    __syncthreads();
#pragma unroll
    for( int i = 0; i < 16; i++, A += lda )
    {
        c[0] += A[0]*bs[i][0]; c[1] += A[0]*bs[i][1]; c[2] += A[0]*bs[i][2]; c[3] += A[0]*bs[i][3];
        c[4] += A[0]*bs[i][4]; c[5] += A[0]*bs[i][5]; c[6] += A[0]*bs[i][6]; c[7] += A[0]*bs[i][7];
        c[8] += A[0]*bs[i][8]; c[9] += A[0]*bs[i][9]; c[10] += A[0]*bs[i][10]; c[11] += A[0]*bs[i][11];
        c[12] += A[0]*bs[i][12]; c[13] += A[0]*bs[i][13]; c[14] += A[0]*bs[i][14]; c[15] += A[0]*bs[i][15];
    } } The bottleneck:
    __syncthreads(); Read A's columns
} while( B != Blast );
for( int i = 0; i < 16; i++, C += ldc )
    C[i] = alpha*c[i] + beta*C[i]; } Do Rank-1 updates
} Store C's block to memory
}

```

Our code vs. CUBLAS 1.1

Performance in multiplying two $N \times N$ matrices on GeForce 8800 GTX:

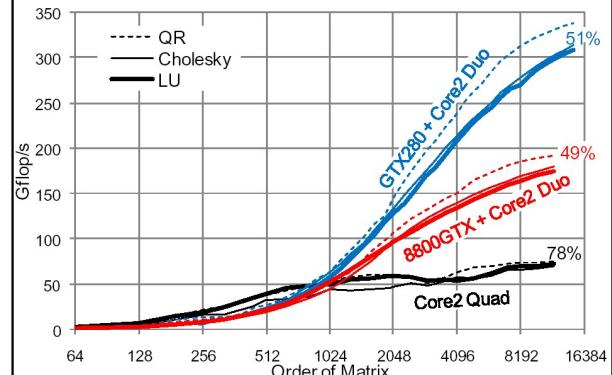


What causes CUBLAS 1.1 to run slower than our code?

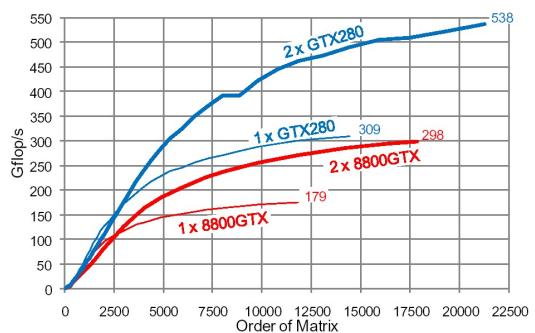
Fast Matrix Factorizations using GPUs

- Use GPU to compute matrix-matrix multiplies only
- Factorize panels on the CPU
- Use look-ahead to overlap computations on CPU and GPU
- Use right-looking algorithms to have more threads in SGEMM
 - Better load balance in the GPU workload, better latency hiding
- Use row-major layout on GPU in LU factorization
 - Requires extra (but fast) matrix transpose for each CPU-GPU transfer
- Substitute triangular solves $LX=B$ with multiply by L^{-1}
 - Provably stable if we do this only when $\|L^{-1}\| < \text{fixed_threshold}$
 - Small pivot growth nearly always assumes small $\|L^{-1}\|$
 - Accuracy of LU assumes small pivot growth anyway
- Use two-level and variable size blocking as finer tuning
 - Thicker blocks impose lower bandwidth requirements in SGEMM
 - Variable size blocking improves CPU/GPU load balance
- Use column-cyclic layout when computing using two GPUs
 - Requires no data exchange between GPUs in pivoting
 - Cyclic layout is used on GPUs only so does not affect panel factorizations

Performance Results



LU Factorization using Two GPUs



- Second GPU allows 1.7x higher rates
- More than half-teraflop using two GPUs

Advanced Tips

- Every round reading 2KB with 256 threads grouped into 2 blocks, each of 128 threads
- At least 2 blocks if using blockwise synchronizations
- Total number of active threads >192 for pipeline efficiency

X. Cui, Y. Chen and H. Mei. Improving performance of matrix multiplication and FFT on GPU.
15th ICPADS, pp 42-48, IEEE Computer Society, 2009.

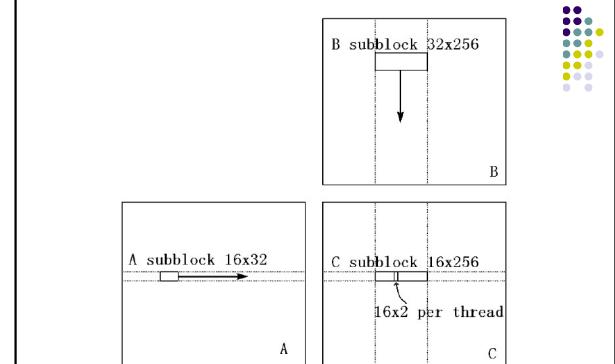


Figure 3: Matrix multiplication: $A \times B = C$ with all arrays in row major

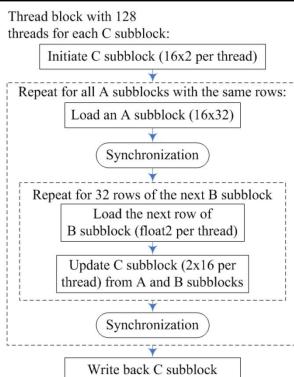
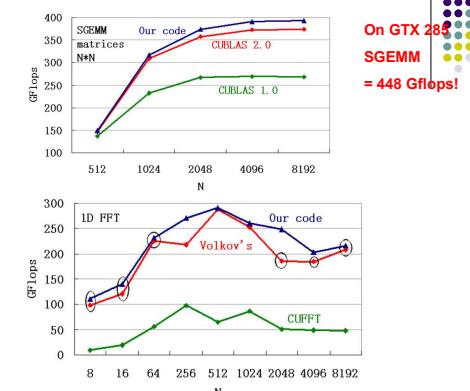
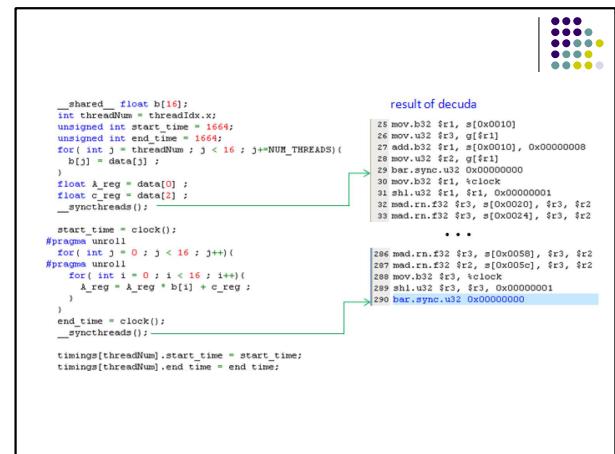
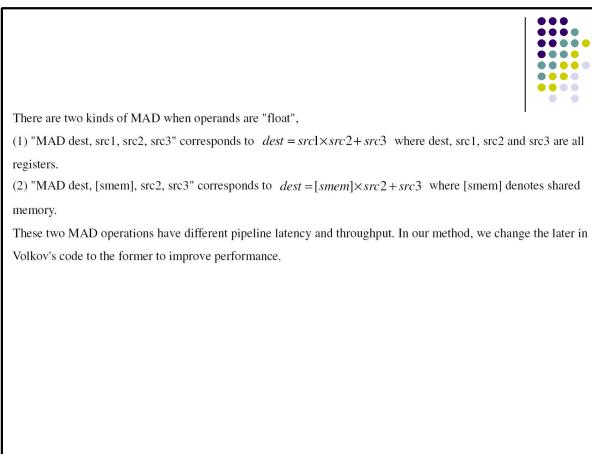


Figure 4: Thread block function for matrix-matrix multiplication.



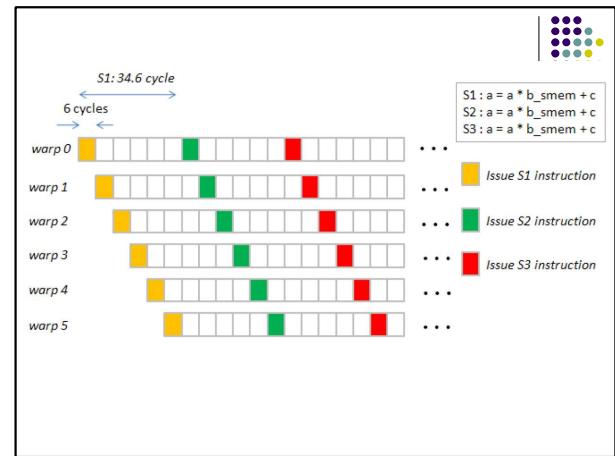
Hand-Tuned SGEMM on GT200 GPU
Lung-Sheng Chien
Department of Mathematics, Tsing Hua university, R.O.C. (Taiwan)

	GTX295 ¹	GTX285	TeslaC1060
# of Streaming Processor	240	240	240
Core Frequency	1242MHz	1476 MHz	1.3 GHz
Memory Speed	999MHz	1242 MHz	800 MHz
Memory Interface	448-bit (7 channel)	512-bit (8 channel)	512-bit (8 channel)
Memory Bandwidth (GB/s)	112	159	102
SP, peak (Gflop/s)	894	1063	933
SP without dual issue	596.2	708.5	624
DP, peak (Gflop/s)	74.5	88.6	78
DRAM (MByte)	896	1024	4096



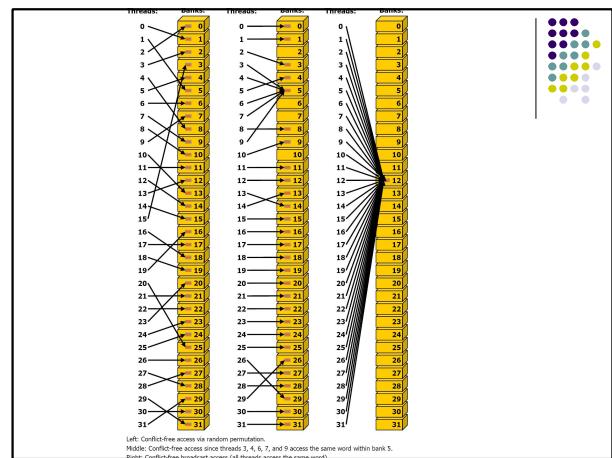
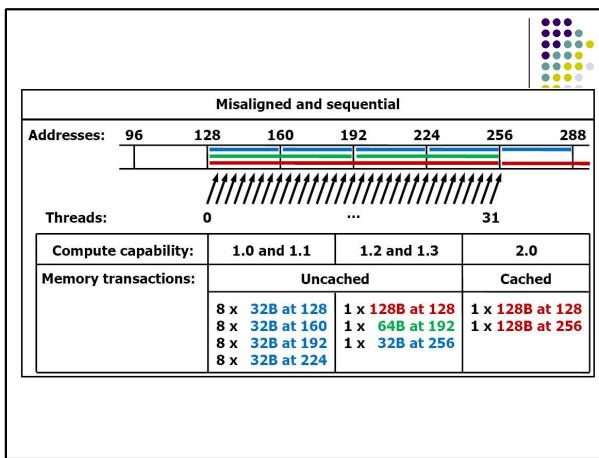
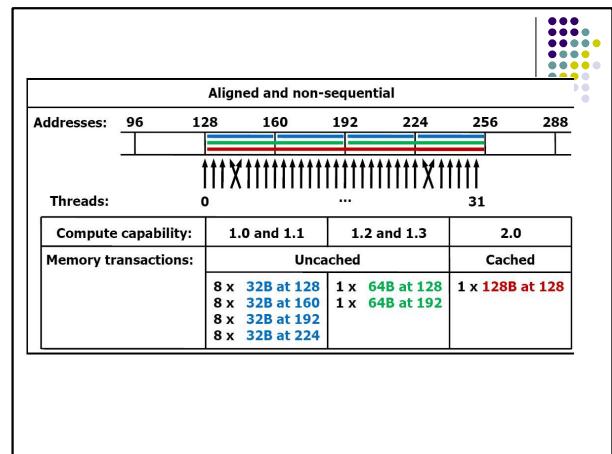
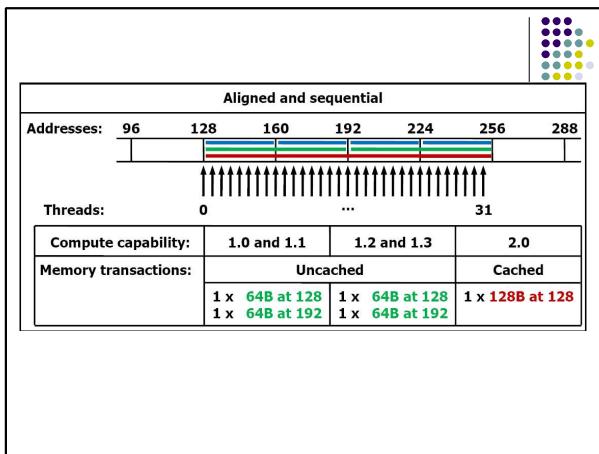
NUM_THREADS	1	64	128	192	224	256	288	320	384	512
Minimum time	34.6	34.6	34.7	38.6	42.4	46.6	54.0	60.2	72.1	96.1
Maximum time	34.6	34.6	34.8	39.3	43.2	49.5	54.7	60.6	72.7	96.9
Total time for one	34.6	34.6	34.6	36	42	48	54	60	72	96
a = a*b_smem + c										

Table 2: average number of cycles per "MAD dest, [smem], src2, src3" on TeslaC1060. Pipeline latency is 34.6 cycle and throughput is 6 cycle /warp.



	Compute Capability	Number of Multiprocessors	Number of CUDA Cores
GeForce GTX 295	1.3	2x30	2x240
GeForce GTX 285, GTX 280	1.3	30	240
GeForce GTX 260	1.3	24	192
GeForce 9800 GX2	1.1	2x16	2x128
GeForce GTS 250, GTS 150, 9800 GTX, 9800 GTX+, 8800 GTS 512	1.1	16	128
GeForce 8800 Ultra, 8800 GTX	1.0	16	128
GeForce 9800 GT, 8800 GT, GTX 280M, 9800M GTX	1.1	14	112
Tesla S1070	1.3	4x30	4x240
Tesla C1060	1.3	30	240

Technical Specifications	Compute Capability			
	1.0	1.1	1.3	2.0
Maximum x- or y-dimension of a grid of thread blocks		65535		
Maximum number of threads per block	512		1024	
Maximum x- or y-dimension of a block	512		1024	
Maximum z-dimension of a block	64			
Warp size	32			
Maximum number of resident blocks per multiprocessor	8			
Maximum number of resident warps per multiprocessor	24	32	48	
Maximum number of resident threads per multiprocessor	768	1024	1536	
Number of 32-bit registers per multiprocessor	8 K	16 K	32 K	
Maximum amount of shared memory per multiprocessor	16 KB	48 KB		
Number of shared memory banks	16	32		
Amount of local memory per thread	16 KB	512 KB		
Constant memory size	64 KB			
Cache working set per multiprocessor for texture memory	8 KB			
Cache working set per multiprocessor for texture memory	Device dependent, between 6 KB and 8 KB			
Maximum width for a 1D texture reference bound to a CUDA array	8192	32768		
Maximum width for a 1D texture reference bound to linear memory	2^{17}			
Maximum width and height for a 2D texture reference bound to linear memory or a CUDA array	65536 x 32768	65536 x 65536		
Maximum width, height, and depth for a 3D texture reference bound to linear memory or a CUDA array	2048 x 2048 x 2048	4096 x 4096 x 4096		
Maximum number of instructions per kernel	2 million			



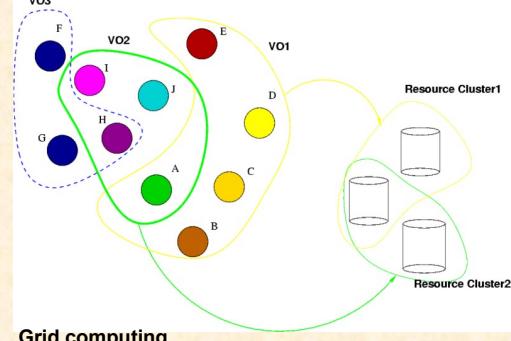
了解各种（非众核） 并行程序设计思想

Programmable Memory Hierarchy

- Registers
- Cache
- Physical memory
- Virtual memory
- Local disk storage
- RAID (via LAN) / Distributed shared-memory
- Data Grid

Ideas of parallel computation

- Multi-threading (e.g. Java)
- Concurrent processes (e.g. Unix)
- SIMD (e.g. Connection Machine-2)
- SPMD (e.g. MPI)
- MIMD (e.g. CORBA)
- Message-passing
- Memory-sharing
- Bulk-Synchronous Parallelism
- Work-flow task parallelism, data parallelism



MPI集群编程示例

IB/verbs 编程示例

```

while (totcount < (user_param->iters * user_param->nmsgops) || totcount < (user_param->iters * user_param->nmsgops)) {
    /* main loop to run over all the qps and post each time n messages */
    for (index = 0; index < user_param->iters * user_param->nmsgops; index++) {
        /* get user param, remap with rem dest index if needed */
        ctx->vc->rma.rkey = rem_dest[index]->rkey;
        qp = ctx->qp[index];
        /* get user param, remap with rem dest index if needed */
        while (ctx->scnt[index] < user_param->iters * (ctx->scnt[index] - ctx->scnt[index]) < user_param->nmsgops) {
            tposted(totcount) = get_cycles();
            if (ibv_post_send(qp, ctx->vc, ibdev_id)) {
                fprintf(stderr, "couldn't post send! qp index = %d qp scnt=%d total scnt %d\n",
                        index, ctx->scnt[index], totcount);
                return 1;
            }
            ctx->scnt[index]++;
            totcount++;
        }
        /* finished posting now polling */
        if (totcount < (user_param->iters * user_param->nmsgops)) {
            do {
                ne = ibv_poll_cq(ctx->cq, 1, &mc);
                /* wait until ne > 0 */
                completed(totcount) = get_cycles(); // MAX wait for the timing
                if (ne < 0) {
                    fprintf(stderr, "poll CQ failed %d\n", ne);
                    return 1;
                }
            } while (ne == 0);
        }
    }
}

```

Pthread多核编程示例

```

• #include <stdio.h>
#include <stdlib.h>
#include <pthread.h>

void *print_message_function( void *ptr );

main()
{
    pthread_t thread1, thread2;
    char *message1 = "Thread 1";
    char *message2 = "Thread 2";
    int iret1, iret2;

    /* Create independent threads each of which will execute function */
    iret1 = pthread_create( &thread1, NULL, print_message_function, (void*) message1);
    iret2 = pthread_create( &thread2, NULL, print_message_function, (void*) message2);

    /* Wait till threads are complete before main continues. Unless we */
    /* wait we run the risk of executing an exit which will terminate */
    /* the process and all threads before the threads have completed. */

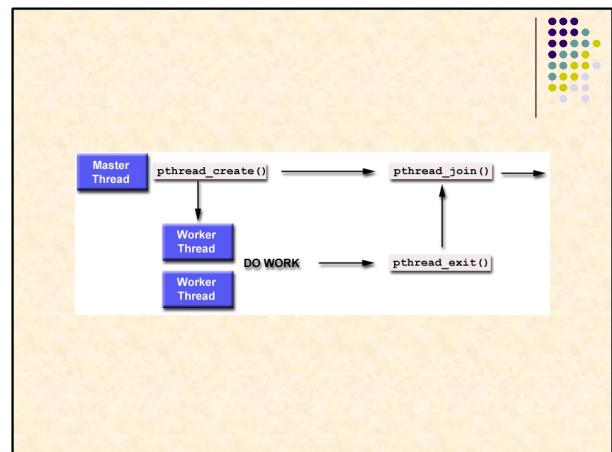
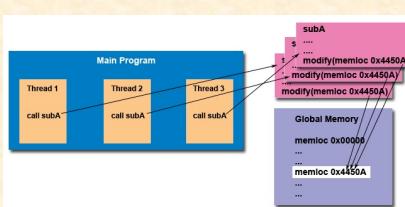
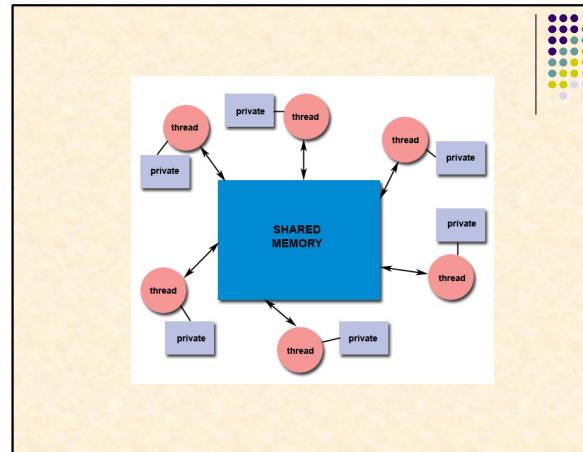
    pthread_join( thread1, NULL);
    pthread_join( thread2, NULL);

    printf("Thread 1 returns: %d\n",iret1);
    printf("Thread 2 returns: %d\n",iret2);
    exit(0);
}

void *print_message_function( void *ptr )
{
    char *message;
    message = (char *)ptr;
    printf("%s \n", message);
}

```

PTHREAD



Thread 1	Thread 2	Balance
Read balance: \$1000		\$1000
	Read balance: \$1000	\$1000
	Deposit \$200	\$1000
Deposit \$200		\$1000
Update balance \$1000+\$200		\$1200
	Update balance \$1000+\$200	\$1200

```
#include <unistd.h> /* Symbolic Constants */
#include <sys/types.h> /* Primitive System Data Types */
#include <pthread.h> /* POSIX Threads */
#include <semaphore.h> /* Semaphore */

void handler (void *ptr);
sem_t mutex;int counter; /* shared variable */

int main(){ int i[2];
pthread_t thread_a; pthread_t thread_b;
i[0] = 0; /* argument to threads */ i[1] = 1;
sem_init(&mutex, 0, 1); /* initialize mutex to 1 - semaphore */
pthread_create (&thread_a, NULL, (void *) &handler, (void *) &i[0]);
pthread_create (&thread_b, NULL, (void *) &handler, (void *) &i[1]);
pthread_join(thread_a, NULL); pthread_join(thread_b, NULL);
sem_destroy(&mutex); /* destroy semaphore */
exit(0);} /* main() */
```

```
void handler (void *ptr) {
    int x; x = *((int *) ptr);
    printf("Thread %d: Waiting to enter critical region...\n", x);
    sem_wait(&mutex); /* down semaphore */
/* START CRITICAL REGION */
    printf("Thread %d: Now in critical region...\n", x);
    printf("Thread %d: Counter Value: %d\n", x, counter);
    printf("Thread %d: Incrementing Counter...\n", x); counter++;
    printf("Thread %d: New Counter Value: %d\n", x, counter);
    printf("Thread %d: Exiting critical region...\n", x);
/* END CRITICAL REGION */
    sem_post(&mutex); /* up semaphore */
    pthread_exit(0); /* exit thread */
}
```

MPI

A Typical MPI Program

```
#include "mpi.h"
...
main(int argc, char** argv) {
    ...
    /* No MPI functions called before this */
    MPI_Init(&argc, &argv);
    ...
    MPI_Finalize();
    /* No MPI functions called after this */
} /* main */
```

```
#include <stdio.h>
#include "mpi.h"
main(int argc, char** argv) {
    int my_rank; /* Rank of process */
    int p; /* Number of processes */
    int source; /* Rank of sender */
    int dest; /* Rank of receiver */
    int tag = 50; /* Tag for messages */
    char message[100]; /* Storage for the message */
    MPI_Status status; /* Return status for receive */
    MPI_Init(&argc, &argv);
    MPI_Comm_rank(MPI_COMM_WORLD, &my_rank);
    MPI_Comm_size(MPI_COMM_WORLD, &p);
    if (my_rank != 0) {
        sprintf(message, "Greetings from process %d!",
               my_rank);
        dest = 0;
        /* Use strlen(message)+1 to include '\0' */
        MPI_Send(message, strlen(message)+1, MPI_CHAR, dest,
                 tag, MPI_COMM_WORLD);
    } else { /* my_rank == 0 */
        for (source = 1; source < p; source++) {
```

```

        MPI_Recv(message, 100, MPI_CHAR, source, tag,
                 MPI_COMM_WORLD, &status);
        printf("%s\n", message);
    }

    MPI_Finalize();
} /* main */

```



When the program is compiled and run with two processes, the output should be

Greetings from process 1!

If it's run with four processes, the output should be

Greetings from process 1!

Greetings from process 2!

Greetings from process 3!

MPI_Send and MPI_Receive

```

int MPI_Send(void* message, int count,
            MPI_Datatype datatype, int dest, int tag,
            MPI_Comm comm)

int MPI_Recv(void* message, int count,
            MPI_Datatype datatype, int source, int tag,
            MPI_Comm comm, MPI_Status* status)

```



MPI datatype	C datatype
MPICHAR	signed char
MPLSHORT	signed short int
MPLINT	signed int
MPLLONG	signed long int
MPLUNSIGNED_CHAR	unsigned char
MPLUNSIGNED_SHORT	unsigned short int
MPLUNSIGNED	unsigned int
MPLUNSIGNED_LONG	unsigned long int
MPLFLOAT	float
MPIDOUBLE	double
MPLLONG_DOUBLE	long double
MPLBYTE	
MPLPACKED	

(Explicit) Synchronisation

```

int MPI_Barrier(MPI_Comm comm)

```

MPI_Barrier provides a mechanism for synchronizing all the processes in the communicator **comm**. Each process blocks (i.e., pauses) until every process in **comm** has called **MPI.Barrier**.



Broadcasting

A communication pattern that involves all the processes in a communicator is a *collective communication*. As a consequence, a collective communication usually involves more than two processes. A *broadcast* is a collective communication in which a single process sends the same data to every process. In MPI the function for broadcasting data is **MPI.Bcast**:

```

int MPI_Bcast(void* message, int count,
              MPI_Datatype datatype, int root, MPI_Comm comm)

```

```

void Get_data2(int my_rank, float* a_ptr, float* b_ptr,
int* n_ptr) {
int root = 0; /* Arguments to MPI_Bcast */
int count = 1;

if (my_rank == 0)
{
    printf("Enter a, b, and n\n");
    scanf("%f %f %d", a_ptr, b_ptr, n_ptr);
}
MPI_Bcast(a_ptr, 1, MPI_FLOAT, root,
MPI_COMM_WORLD);
MPI_Bcast(b_ptr, 1, MPI_FLOAT, root,
MPI_COMM_WORLD);
MPI_Bcast(n_ptr, 1, MPI_INT, root,
MPI_COMM_WORLD);
} /* Get_data2 */

```

Reduction

```

int MPI_Reduce(void* operand, void* result,
int count, MPI_Datatype datatype, MPI_Op op,
int root, MPI_Comm comm)

```

Operation Name	Meaning
MPI_MAX	Maximum
MPI_MIN	Minimum
MPI_SUM	Sum
MPI_PROD	Product
MPI_LAND	Logical And
MPI_BAND	Bitwise And
MPI_LOR	Logical Or
MPI_BOR	Bitwise Or
MPI_LXOR	Logical Exclusive Or
MPI_BXOR	Bitwise Exclusive Or
MPI_MAXLOC	Maximum and Location of Maximum
MPI_MINLOC	Minimum and Location of Minimum

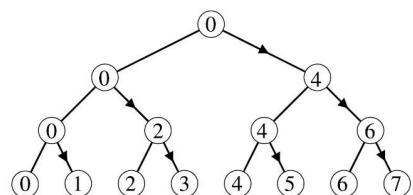
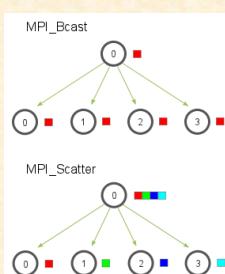


Figure 1: Processors configured as a tree

Yet Another Piece of Advice

Advice to implementors. It is strongly recommended that MPI_REDUCE be implemented so that the same result be obtained whenever the function is applied on the same arguments, appearing in the same order. Note that this may prevent optimizations that take advantage of the physical location of processors. (End of advice to implementors.)

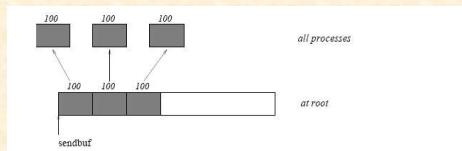


Scattering

```

int MPI_Scatterv(void* sendbuf, int *sendcounts, int *displs,
MPI_Datatype sendtype, void* recvbuf, int recvcount,
MPI_Datatype recvtype, int root, MPI_Comm comm)

```



The root process scatters sets of 100 ints to each process in the group.

Gathering

```
int MPI_Gatherv(void* sendbuf, int sendcount, MPI_Datatype sendtype,
                void* recvbuf, int *recvcounts, int *displs,
                MPI_Datatype recvtype, int root, MPI_Comm comm)
```

The root process gathers 100 ints from each process in the group.

MPI_Allgather

```
MPI_Allgather( void* send_data,
                int send_count,
                MPI_Datatype send_datatype,
                void* recv_data,
                int recv_count,
                MPI_Datatype recv_datatype,
                MPI_Comm communicator
            )
```

MPI_Allgather / MPI_Alltoall

```
MPI_Allgather(a,m,1,MPI_INT,b,1,MPI_INT,MPI_COMM_WORLD);
MPI_Alltoall(a,1,MPI_INT,b,1,MPI_INT,MPI_COMM_WORLD);
```

MPI 实现 alltoall.c

```

520    /* Do the pairwise exchanges */
521    for (i=1; i<comm_size; i++) {
522        if (poff == 1) {
523            /* use exclusive-or algorithm */
524            src = dat = rank ^ i;
525        } else {
526            src = (rank - i + comm_size) % comm_size;
527            dat = (rank + i) % comm_size;
528        }
529
530        mpi_errno = MPI_C_Sendrecv_ft(((char *)sendbuf +
531                                         dst*sendcount*sendtype extent),
532                                         sendcount, sendtype, dst,
533                                         MPIR_ALLTOALL_TAG,
534                                         ((char *)recvbuf +
535                                         src*recvcount*recvtype extent),
536                                         recvcount, recvtype, src,
537                                         MPIR_ALLTOALL_TAG, comm, &status, errflag);
538
539        if (mpi_errno) {
540            /* for communication errors, just record the error but continue */
541            *errflag = TRUE;
542            MPIU_ERR_SET(mpi_errno, MPI_ERR_OTHER, "**fail");
543            MPIU_ERR_ADD(mpi_errno_ret, mpi_errno);
544        }
545    }

```

Communication Groups

```

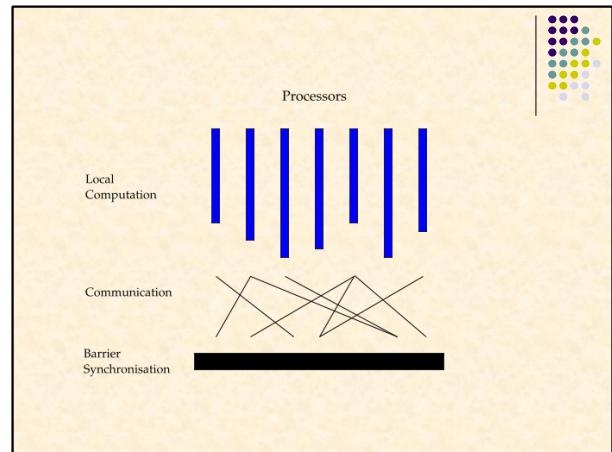
MPI_Group MPI_GROUP_WORLD;
MPI_Group first_row_group;
MPI_Comm first_row_comm;
int row_size;
int* process_ranks;

/* Make a list of the processes in the new
 * communicator */
process_ranks = (int*) malloc(q*sizeof(int));
for (proc = 0; proc < q; proc++)
    process_ranks[proc] = proc;

/* Get the group underlying MPI_COMM_WORLD */
MPI_Comm_group(MPI_COMM_WORLD, &MPI_GROUP_WORLD);

/* Create the new group */
MPI_Group_incl(MPI_GROUP_WORLD, q, process_ranks, &first_row_group);
MPI_Comm_create(MPI_COMM_WORLD, first_row_group,
               &first_row_comm);

```



GPU集群LINPACK

Hardware Configuration

1. SUN Ultra 24 workstation with an Intel Core2 Extreme Q6850 (3.0GHz) CPU and 8GB of memory plus a Tesla C1060 card.
2. Cluster with 8 nodes, each node connected to half of a Tesla S1070 system, containing 4 GPUs, so that each node is connected to 2 GPUs. Each node has 2 Intel Xeon E5462 (2.8GHz with 1600MHz FSB) and 16GB of memory. The nodes are connected with SDR (Single Data Rate) Infiniband.

PCI Bandwidth

Sun Ultra 24		
	Pageable memory	Pinned memory
H2D	2132 MB/s	5212 MB/s
D2H	1882 MB/s	5471 MB/s

Supermicro 6015TW		
	Pageable memory	Pinned memory
H2D	2524 MB/s	5651 MB/s
D2H	2084 MB/s	5301 MB/s

LINPACK (wiki)

- The **LINPACK Benchmarks** are a measure of a system's floating point computing power. Introduced by Jack Dongarra, they measure how fast a computer solves a dense N by N system of linear equations $Ax = b$, which is a common task in engineering. The solution is obtained by Gaussian elimination with partial pivoting, with $2/3 \cdot N^3 + 2 \cdot N^2$ floating point operations. The result is reported in millions of floating point operations per second.
- Massimiliano Fatica, *Accelerating linpack with CUDA on heterogenous clusters*, GPGPU'09.
Slides extracted from the above paper

Pinned Memory

```
// Regular malloc/free
double *A;
A = malloc(N*N*sizeof(double));
free(A);

// Page-locked version
double *A;
cudaMallocHost(A,N*N*sizeof(double));
cudaFreeHost(A);
```

Solving Linear Equations

Figure 1: LU factorization: the grey area represents the portion of the matrix already factored. The red area is the current block being factorized. Once this factorization is ready, it is applied to the sub-matrix on the right. The final step is to update the trailing sub-matrix in yellow.

LU Decomposition

$L \cup U = A$

$Ax = (L \cup U)x = L(Ux) = b$

$$\begin{bmatrix} l_{11} & 0 & 0 \\ l_{21} & l_{22} & 0 \\ l_{31} & l_{32} & l_{33} \end{bmatrix} \begin{bmatrix} u_{11} & u_{12} & u_{13} \\ 0 & u_{22} & u_{23} \\ 0 & 0 & u_{33} \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix}$$

$$\begin{bmatrix} l_{11} u_{11} & l_{11} u_{12} & l_{11} u_{13} \\ l_{21} u_{11} + l_{22} u_{22} & l_{21} u_{12} + l_{22} u_{22} & l_{21} u_{13} + l_{22} u_{23} \\ l_{31} u_{11} + l_{32} u_{22} & l_{31} u_{12} + l_{32} u_{22} + l_{33} u_{33} & l_{31} u_{13} + l_{32} u_{23} + l_{33} u_{33} \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix}$$

first solve $Ly = b$ for y . This can be done by forward substitution

$$y_1 = \frac{b_1}{l_{11}}$$

$$y_i = \frac{1}{l_{ii}} \left(b_i - \sum_{j=1}^{i-1} l_{ij} y_j \right)$$

for $i = 2, \dots, N$. Then solve $Ux = y$ for x . This can be done by back substitution

$$x_N = \frac{y_N}{u_{NN}}$$

$$x_i = \frac{1}{u_{ii}} \left(y_i - \sum_{j=i+1}^N u_{ij} y_j \right)$$

for $i = N-1, \dots, 1$.

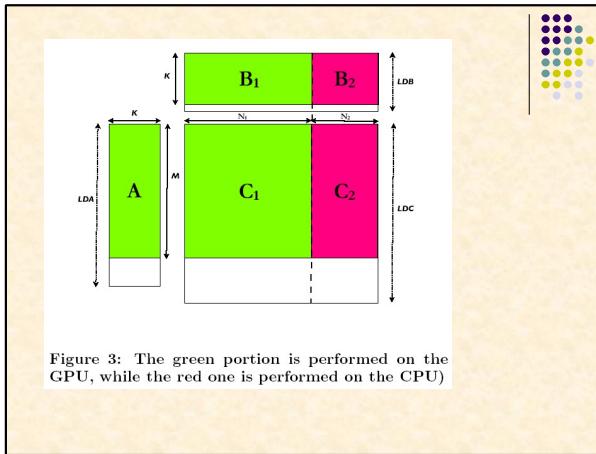


Figure 3: The green portion is performed on the GPU, while the red one is performed on the CPU)

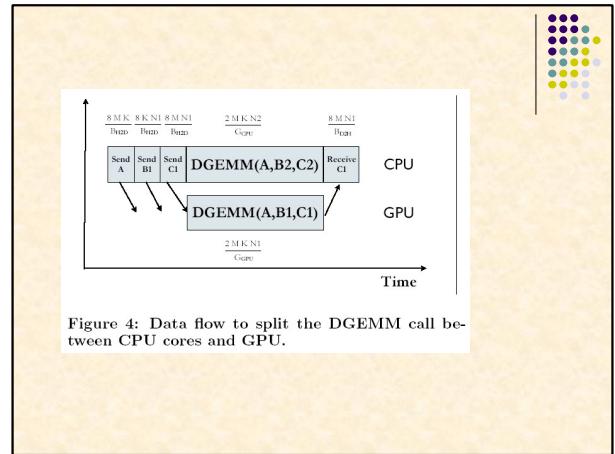


Figure 4: Data flow to split the DGEMM call between CPU cores and GPU.

B_{H2D} : Bandwidth from host to device expressed in GB/s

G_{GPU} : Sustained performance of DGEMM on the GPU expressed in $GFlops$

G_{CPU} : Sustained performance of DGEMM on the CPU expressed in $GFlops$

B_{D2H} : Bandwidth from device to host expressed in GB/s

A DGEMM call on the host CPU performs $2KMN$ operations, so if the CPU cores can perform this operation at G_{CPU} the total time is:

$$T_{CPU}(M, K, N) = 2 \frac{MKN}{G_{CPU}}$$

The total time to offload a DGEMM call to the GPU has an I/O component that accounts for both the data transfer from the CPU memory space to the GPU memory space and vice versa plus a computational part once the data is on the GPU. We can express this time as:

$$T_{GPU}(M, K, N) = 8 \frac{(MK + KN + MN)}{B_{H2D}} + 2 \frac{MKN}{G_{GPU}} + 8 \frac{(MN)}{B_{D2H}}$$

the factor 8 is the size of a double in bytes.
The optimal split will be

$$T_{CPU}(M, K, N2) = T_{GPU}(M, K, N1) \quad \text{with} \quad N = N1 + N2$$

For an initial approximation of the optimal split fraction $\eta = N1/N$, we can omit the transfer time ($O(N^2)$) compared to the computation ($O(N^3)$). From a simple manipulation, the optimal split is

$$\eta = \frac{G_{GPU}}{G_{GPU} + G_{CPU}}$$

On the cluster, where the quad core Xeon has a DGEMM performance of 40 GFlops and the GPU a DGEMM performance of 82 GFlops, this formula predicts $\eta = 0.67$, very close to the optimal value of 0.68 found by experiments.

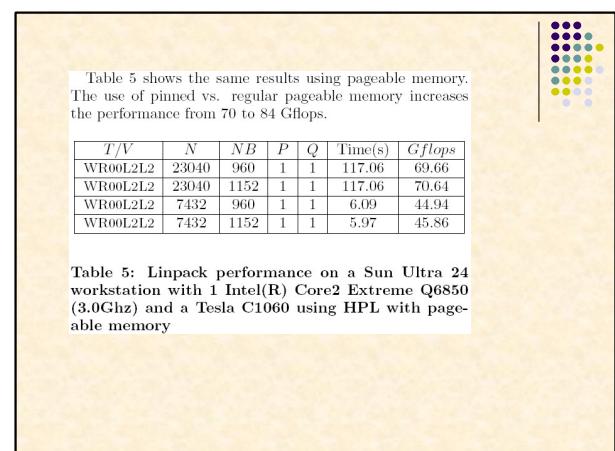
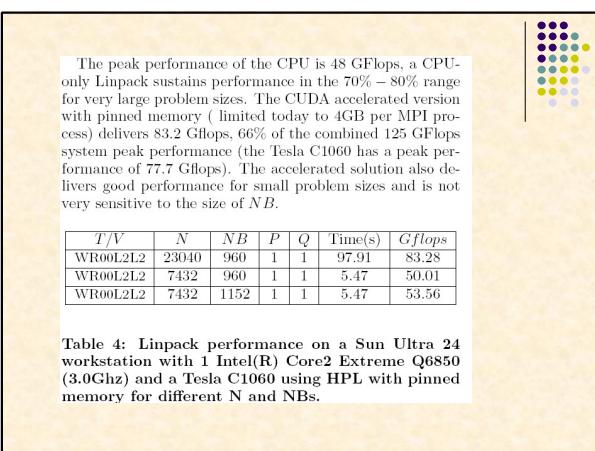
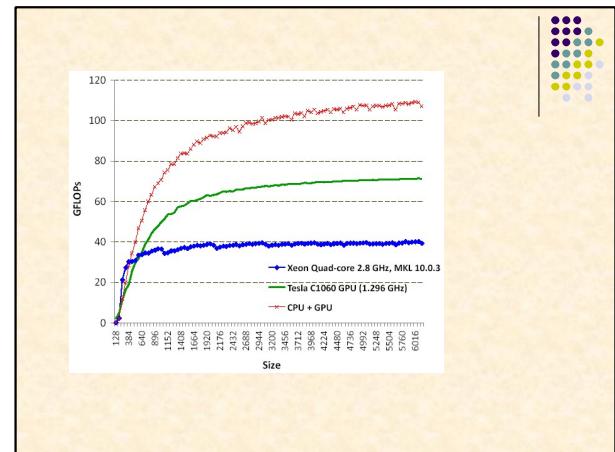
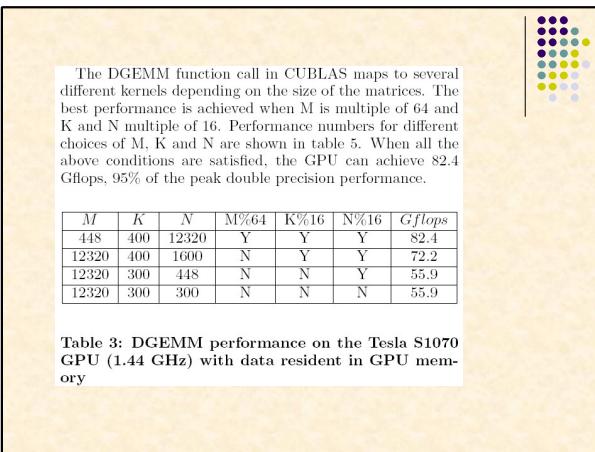
```
// Copy A from CPU memory to GPU memory devA
status = cublasSetMatrix (m, k , sizeof(A[0]), A, lda, devA, m_gpu);
// Copy B1 from CPU memory to GPU memory devB
status = cublasSetMatrix (k ,n_gpu, sizeof(B[0]), B, ldb, devB, k_gpu);
// Copy C1 from CPU memory to GPU memory devC
status = cublasSetMatrix (m, n_gpu, sizeof(C[0]), C, ldc, devC, m_gpu);

// Perform DGEMM(devA,devB,devC) on GPU
// Control immediately return to CPU
cublasDgemm('n', 'n', m, n_gpu, k, alpha, devA, m, devB, k, beta, devC, m);

// Perform DGEMM(A,B2,C2) on CPU
dgemm_cpu('n','n',m,n_cpt,k, alpha, A, lda,B1db*n_gpu, ldb, beta,C+ldc*n_gpu, ldc);

// Copy devC from GPU memory to CPU memory C1
status = cublasGetMatrix (m, n, sizeof(C[0]), devC, m, C, *ldc);
```

It turns out that on Intel systems using Front Side Bus (FSB), it is better not to overlap the transfer to the GPU with computations on the CPU (the memory system cannot supply data to both the PCIe and the CPU at maximum speed).



On the workstation, the biggest problem that can be solved with the available memory is $N = 32320$ and the Linpack score is now 90 Gflops, 72% of peak performance.

Sun Ultra 24						
T/V	N	NB	P	Q	Time(s)	Gflops
WR00R2L2	32320	1152	1	1	250.01	90
Cluster						
T/V	N	NB	P	Q	Time(s)	Gflops
WR11R2L2	118114	960	4	4	874	1258

Table 8: Linpack performance using HPL with pinned memory (pre-release CUDA version)



一致性 时态逻辑 安全性与进展性



用一条直线来表示系统中每一个处理器的运行，时间从左边开始到右边。每一个共享内存的操作我们把它写在处理器的直线上。两个主要的操作为“读”和“写”，用下面的表达式表示：

W(var) value

该表达式的含义为：将值value写到共享变量var中，而

R(var) value

该表达式的含义为：读取共享变量var，获取它的值value。比如：W(x)1 表示：将1写到x中，而R(y)3表示：读取变量y的值3。

更多的操作（特别是同步操作）在需要的时候我们再来定义他的记号。为简单起见，假设所有的变量都初始化为0。我们要特别注意一点，在高级语言中的一条语句（比如 $x = x + 1;$ ）通常将会涉及到几次内存操作。如果x之前为0，则那条高级语言中的语句将变成（不考虑其它处理器）：

P1: R(x)0 W(x)1

严格一致性(Strict Consistency):



- 任意read(x)操作都要读到最新的write(x)的结果。依赖于绝对的全局时钟。

P1: W(x)1

P2: R(x)1 R(x)1

这表示，“处理器P1将1写到变量x中；一段时间之后处理器P2读取x的值1。然后再读取一次，获得相同的值。

我们再给出一个符合“内存严格一致性模型”的场景：

P1: W(x)1

P2: R(x)0 R(x)1

这一次，处理器P2先执行，它先读取x的值0，当它第二次读取x时却获得了处理器P1写入的x的值1。注意这两个场景能够通过将同一个程序在同一个处理器上执行两次而获得。

我们给出一个不符合“内存严格一致性模型”的场景：

P1: W(x)1

P2: R(x)0 R(x)1

顺序一致性 (Sequential Consistency):



- “（并发程序在多处理器上的）任何一次执行结果都相同，就像所有处理器的操作按照某个顺序执行，各个微处理器的操作按照其程序指定的顺序进行。”。

在一个写操作发生之前看到该写操作的结果是可能的，

P1: W(x)1

P2: R(x)1

P1: W(x)1

P2: R(x)1 R(x)2

P3: R(x)1 R(x)2

P4: W(x)2

这个场景是合法的顺序一致性内存模型的原因是，下面的交替操作在严格一致性内存模型中将会是合法的：

P1: W(x)1

P2: R(x)1 R(x)2

P3: R(x)1 R(x)2

P4: W(x)2

下面是一个不符合顺序一致性内存模型的场景：

P1: W(x)1

P2: R(x)1 R(x)2

P3: R(x)2 R(x)1

P4: W(x)2

顺序一致性在分布式算法中的应用例

- 最短路径迭代计算收敛终止判定
 - 每个线程独立反复运行Dijkstra算法，并记录第k轮有多少结点被更新；
 - 无更新且其他线程见过该线程最新一轮则进入终止准备状态；
 - 一个线程见到其他线程相关结点更新后先退出终止准备状态，再告知对方已经见过。
 - 原子性判定所有线程都终止准备状态则终止。

弱一致性 (Weak Consistency):

- 如果我们仅仅将竞争访问分为同步与非同步访问，并且同时要求符合下列条件，那么我们就得到了“弱一致性”：
 - 对同步变量的访问是“顺序一致性”的
 - 直到之前对所有同步变量的写操作完成之后，才允许访问这些同步变量。
 - 直到之前对同步变量的访问完成之后，才允许我们访问（读或写）这些同步变量。

下面给出一个符合“弱一致性”的场景，显示了“弱一致性”的真正用途：

P1: W(x)1 W(x)2

P2: R(x)0 R(x)2 S R(x)2

P3: R(x)1 S R(x)2

换一句话说，根本没有要求一个处理器广播它对变量的修改，直到一个同步访问的发生。在一个基于网络而不是总线的分布式系统中，这能够极大的减少信息的互通（注意到，在现实中没有人会故意写一个具有这样行为的程序；你绝对不想去读一个别人正在更新的变量。读操作必须发生在S之后。我提到过一些同步算法，比如松弛算法，它不要求内存一致性的概念。这些算法在“弱一致性”系统中不能工作，因为在“弱一致性”系统中推迟了数据的交流直到同步点。）

因果一致性 (Casual Consistency)

- 有因果关系的写操作必须按照它们的因果关系的顺序被看到，没有因果关系的写操作可以以任意顺序被别的进程看到。

思考题：

- 画出一个满足顺序一致但不满足因果一致的情形
- 最短路径迭代计算收敛终止判定算法可否使用因果一致性

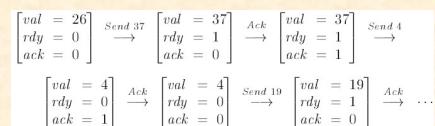
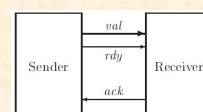
Specifying a Simple Clock

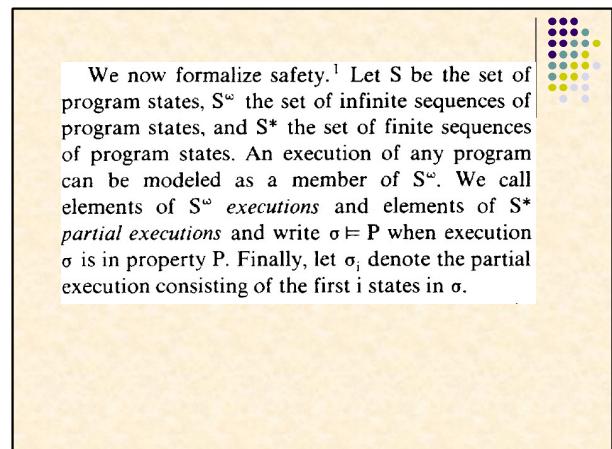
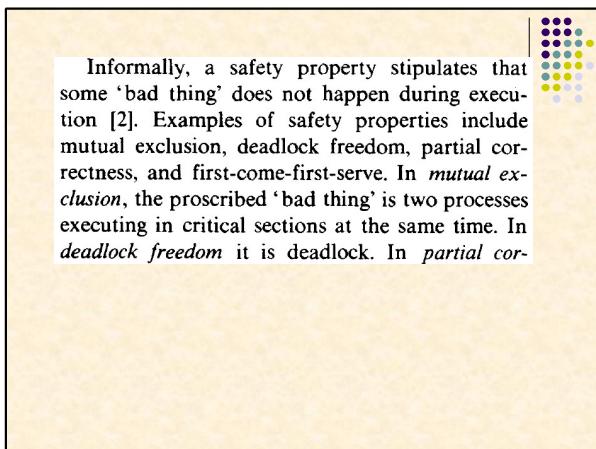
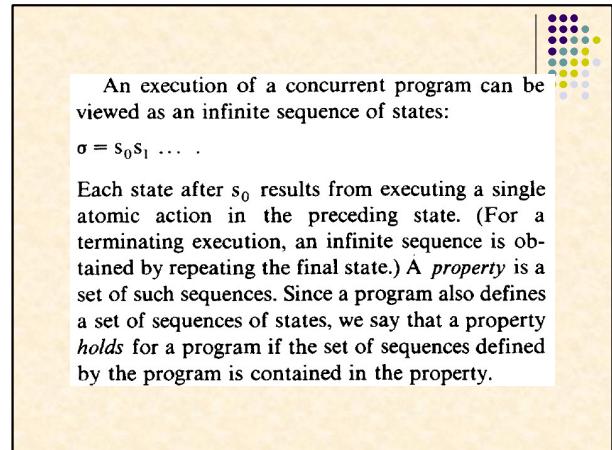
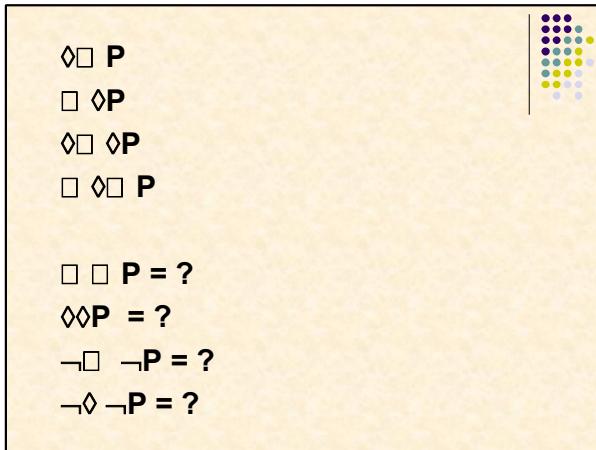
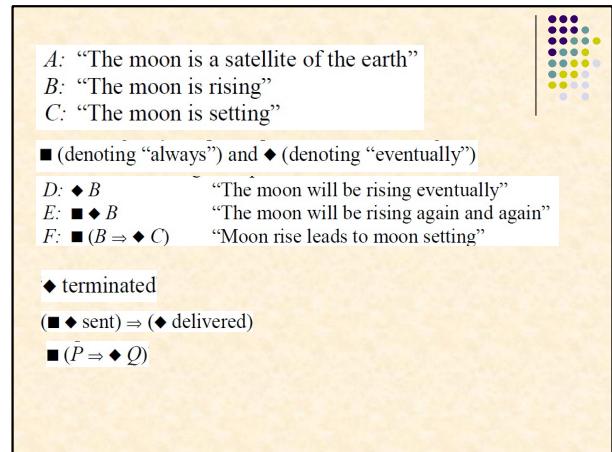
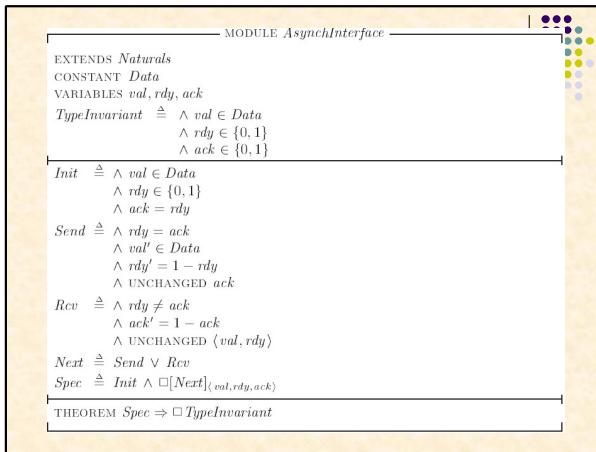
$$HCini \triangleq hr \in \{1, \dots, 12\}$$

$$HCnxt \triangleq hr' = \text{IF } hr \neq 12 \text{ THEN } hr + 1 \text{ ELSE } 1$$

```
MODULE HourClock
EXTENDS Naturals
VARIABLE hr
HCini triangleq hr in {1 .. 12}
HCnxt triangleq hr' = IF hr neq 12 THEN hr + 1 ELSE 1
HC triangleq HCini and [HCnxt]_hr
THEOREM HC implies HCini
```

An Asynchronous Interface





Safety properties



For P to be a safety property, if P does not hold for an execution, then at some point some ‘bad thing’ must happen. Such a ‘bad thing’ must be irremediable because a safety property states that the ‘bad thing’ never happens during execution. Thus, P is a *safety property* if and only if the following holds.

Safety

$$(\forall \sigma : \sigma \in S^\omega : \sigma \not\models P \Rightarrow (\exists i : 0 \leq i : (\forall \beta : \beta \in S^\omega : \sigma, \beta \not\models P))).$$

任何无限运行序列，如若违反P性质，则存在其有限前缀，使得该前缀之所有无限扩展均违反P。

Liveness properties



Informally, a liveness property stipulates that a ‘good thing’ happens during execution [2]. Examples of liveness properties include starvation freedom, termination, and guaranteed service. In *starvation freedom*, which states that a process makes progress infinitely often, the ‘good thing’ is making progress. In *termination*, which asserts that a program does not run forever, the ‘good thing’ is completion of the final instruction. Finally, in *guaranteed service*,² which states that every request for service is satisfied eventually, the ‘good thing’ is receiving service.



We now formalize liveness. A partial execution α is *live* for a property P if and only if there is a sequence of states β such that $\alpha\beta \models P$. A *liveness property* is one for which every partial execution is live. Thus, P is a liveness property if and only if the following holds.

Liveness

$$(\forall \alpha : \alpha \in S^* : (\exists \beta : \beta \in S^\omega : \alpha\beta \models P)).$$

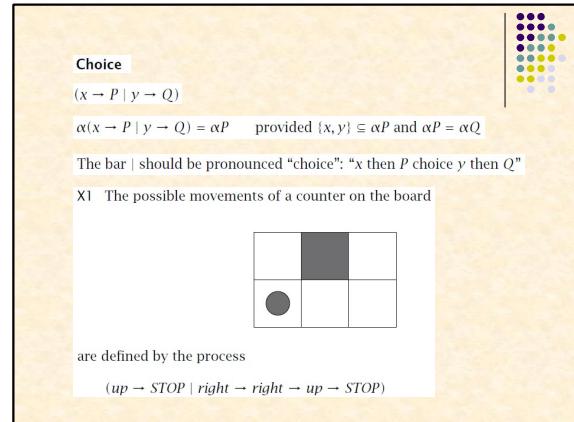
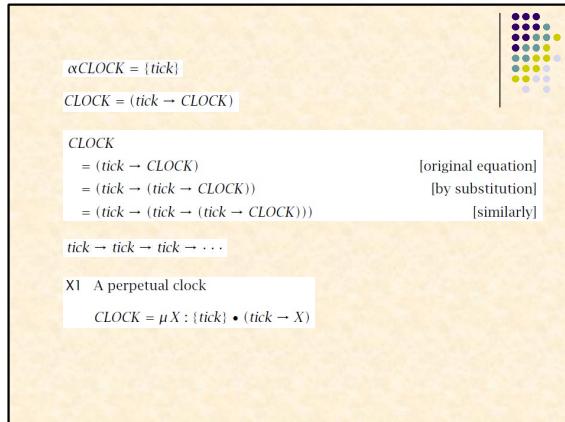
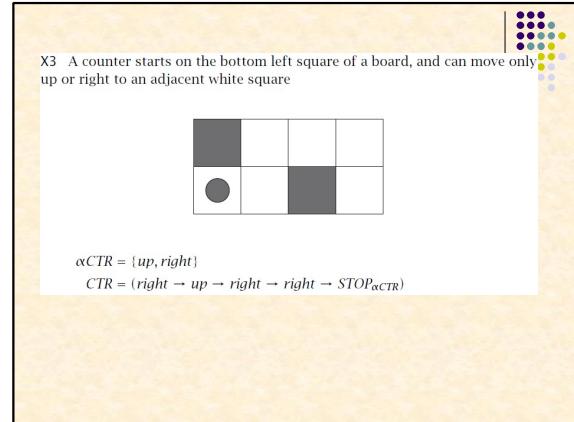
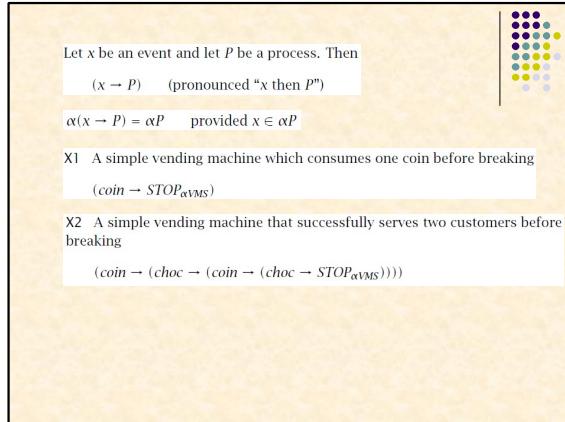
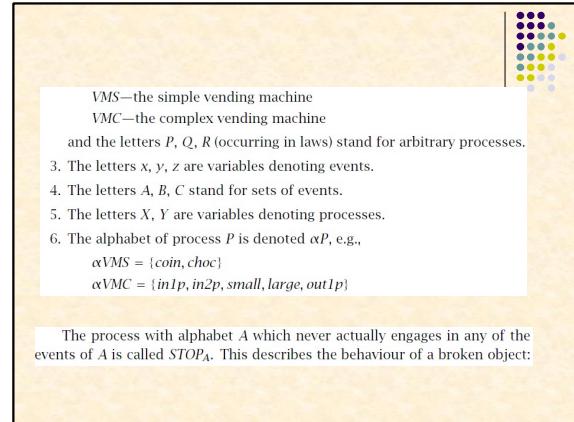
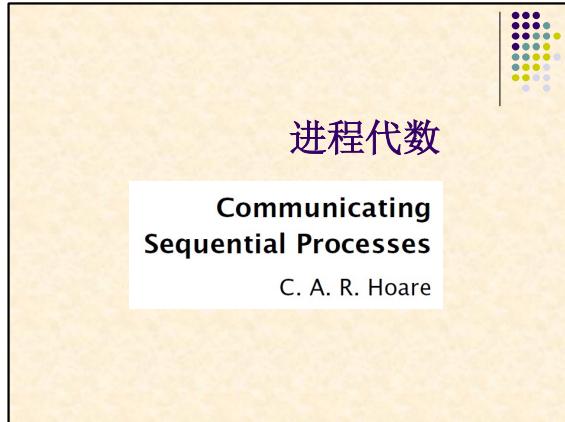
任何有限运行序列均可以被无限扩展为一个满足P的无限序列。

Theorem 1. Every property P is the intersection of a safety property and a liveness property.

Proof. Let \bar{P} be the smallest safety property containing P and let L be $\neg(\bar{P} - P)$. Then,

$$\begin{aligned} L \cap \bar{P} &= \neg(\bar{P} - P) \cap \bar{P} \\ &= (\neg\bar{P} \cup P) \cap \bar{P} \\ &= (\neg\bar{P} \cap \bar{P}) \cup (P \cap \bar{P}) \\ &= P \cap \bar{P} = P. \end{aligned}$$

It remains to show that L is dense, and hence a liveness property. By way of contradiction, suppose there is a nonempty open set O contained in $\neg L$ and thus L is not dense. Then, $O \subseteq (\bar{P} - P)$. Consequently, $P \subseteq (\bar{P} - O)$. The intersection of two closed sets is closed, so $\bar{P} - O$ is closed and thus a safety property. This contradicts the hypothesis that \bar{P} is the smallest safety property containing P . \square



X2 A machine which offers a choice of two combinations of change for 5p (compare 1.1.2 X3 and X4, which offer no choice).

$$CH5C = \mu X \bullet (in.5p \rightarrow (out.1p \rightarrow out.1p \rightarrow out.2p \rightarrow out.2p \rightarrow CH5C) \mid out.2p \rightarrow out.1p \rightarrow out.2p \rightarrow CH5C)$$

The choice is exercised by the customer of the machine. □

X3 A machine that serves either chocolate or toffee on each transaction

$$VMCT = \mu X \bullet coin \rightarrow (choc \rightarrow X \mid toffee \rightarrow X)$$

□

X4 A more complicated vending machine, which offers a choice of coins and a choice of goods and change

$$VMC = (in.2p \rightarrow (large \rightarrow VMC) \mid small \rightarrow out.1p \rightarrow VMC) \mid in.1p \rightarrow (small \rightarrow VMC) \mid in.1p \rightarrow (large \rightarrow VMC) \mid in.1p \rightarrow STOP))$$

X7 A copying process engages in the following events

- in.0*—input of zero on its input channel
- in.1*—input of one on its input channel
- out.0*—output of zero on its output channel
- out.1*—output of one on its output channel

Its behaviour consists of a repetition of pairs of events. On each cycle, it inputs a bit and outputs the same bit

$$COPYBIT = \mu X \bullet (in.0 \rightarrow out.0 \rightarrow X \mid in.1 \rightarrow out.1 \rightarrow X)$$

Pictures

($coin \rightarrow STOP_{\alpha_{VMS}}$)
 $(coin \rightarrow (choc \rightarrow (coin \rightarrow (choc \rightarrow STOP_{\alpha_{VMS}}))))$
 $VMCT = \mu X \bullet coin \rightarrow (choc \rightarrow X \mid toffee \rightarrow X)$

X5 A machine that allows its customer to sample a chocolate, and trusts him to pay after. The normal sequence of events is also allowed

$$VMCRED = \mu X \bullet (coin \rightarrow choc \rightarrow X \mid choc \rightarrow coin \rightarrow X)$$

Laws

L1 $(x : A \rightarrow P(x)) = (y : B \rightarrow Q(y)) \equiv (A = B \wedge \forall x : A \bullet P(x) = Q(x))$

$$(x \rightarrow P \mid y \rightarrow Q) = (y \rightarrow Q \mid x \rightarrow P)$$

$(x \rightarrow P) \neq STOP$

$$(c \rightarrow P) \neq (d \rightarrow Q) \quad \text{if } c \neq d$$

$$(c \rightarrow P) = (c \rightarrow Q) \equiv P = Q$$

$$(coin \rightarrow choc \rightarrow coin \rightarrow choc \rightarrow STOP) \neq (coin \rightarrow STOP)$$

$$\mu X \bullet F(X) = F(\mu X \bullet F(X))$$

Traces

$\langle x, y \rangle$ consists of two events, x followed by y .
 $\langle x \rangle$ is a sequence containing only the event x .
 $\langle \rangle$ is the empty sequence containing no events.

$\langle coin, choc, coin, choc \rangle$
 $s \cap t$
 $\langle coin, choc \rangle \cap \langle coin, toffee \rangle = \langle coin, choc, coin, toffee \rangle$
 $\langle in.1p \rangle \cap \langle in.1p \rangle = \langle in.1p, in.1p \rangle$
 $\langle in.1p, in.1p \rangle \cap \langle \rangle = \langle in.1p, in.1p \rangle$

L1 $s \cap \langle \rangle = \langle \rangle \cap s = s$
L2 $s \cap (t \cap y) = (s \cap t) \cap u$

X1 The only trace of the behaviour of the process $STOP$ is $\langle \rangle$. The notebook of the observer of this process remains forever blank

$$traces(STOP) = \{ \langle \rangle \}$$

□

X2 There are only two traces of the machine that ingests a coin before breaking

$$traces(coin \rightarrow STOP) = \{ \langle \rangle, \langle coin \rangle \}$$

□

X3 A clock that does nothing but $tick$

$$traces(\mu X \bullet tick \rightarrow X) = \{ \langle \rangle, \langle tick \rangle, \langle tick, tick \rangle, \dots \} = \{ tick \}^*$$

As with most interesting processes, the set of traces is infinite, although of course each individual trace is finite.

□

X4 A simple vending machine

$$traces(\mu X \bullet coin \rightarrow choc \rightarrow X) = \{ s \mid \exists n \bullet s \leq (coin, choc)^n \}$$

□

L1 $traces(STOP) = \{ t \mid t = \langle \rangle \} = \{ \langle \rangle \}$

L2 $traces(c \rightarrow P) = \{ t \mid t = \langle \rangle \vee (t_0 = c \wedge t' \in traces(P)) \} = \{ \langle \rangle \} \cup \{ \langle c \rangle \cap t \mid t \in traces(P) \}$

L3 $traces(c \rightarrow P \mid d \rightarrow Q) = \{ t \mid t = \langle \rangle \vee (t_0 = c \wedge t' \in traces(P)) \vee (t_0 = d \wedge t' \in traces(Q)) \}$

If P and Q are processes with the same alphabet, we introduce the notation $P \parallel Q$

to denote the process which behaves like the system composed of processes P and Q interacting in lock-step synchronisation as described above.

Examples

X1 A greedy customer of a vending machine is perfectly happy to obtain a toffee or even a chocolate without paying. However, if thwarted in these desires, he is reluctantly prepared to pay a coin, but then he insists on taking a chocolate

$$GRCUST = (toffee \rightarrow GRCUST \mid choc \rightarrow GRCUST \mid coin \rightarrow choc \rightarrow GRCUST)$$

When this customer is brought together with the machine $VMCT$ (1.1.3 X3) his greed is frustrated, since the vending machine does not allow goods to be extracted before payment. Similarly, $VMCT$ never gives a toffee, because the customer never wants one after he has paid

$$(GRCUST \parallel VMCT) = \mu X \bullet (coin \rightarrow choc \rightarrow X)$$

X2 A foolish customer wants a large biscuit, so he puts his coin in the vending machine VMC . He does not notice whether he has inserted a large coin or a small one; nevertheless, he is determined on a large biscuit

$$FOOLCUST = (in2p \rightarrow large \rightarrow FOOLCUST \mid in1p \rightarrow large \rightarrow FOOLCUST)$$

Unfortunately, the vending machine is not prepared to yield a large biscuit for only a small coin

$$(FOOLCUST \parallel VMC) = \mu X \bullet (in2p \rightarrow large \rightarrow X \mid in1p \rightarrow STOP)$$

L1 $P \parallel Q = Q \parallel P$

The next law shows that when three processes are assembled, it does not matter in which order they are put together

L2 $P \parallel (Q \parallel R) = (P \parallel Q) \parallel R$

Thirdly, a deadlocked process infects the whole system with deadlock; but composition with $RUN_{\alpha P}$ (1.1.3 X8) makes no difference

L3A $P \parallel STOP_{\alpha P} = STOP_{\alpha P}$

L3B $P \parallel RUN_{\alpha P} = P$

The next laws show how a pair of processes either engage simultaneously in the same action, or deadlock if they disagree on what the first action should be

L4A $(c \rightarrow P) \parallel (c \rightarrow Q) = (c \rightarrow (P \parallel Q))$

L4B $(c \rightarrow P) \parallel (d \rightarrow Q) = STOP \quad \text{if } c \neq d$

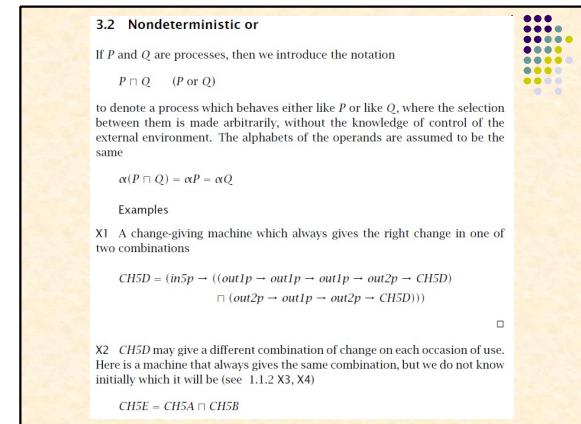
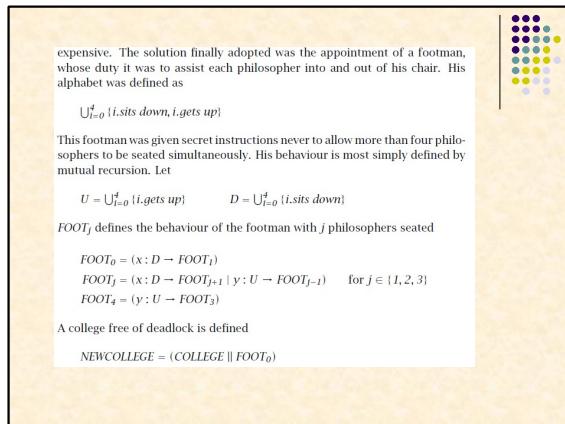
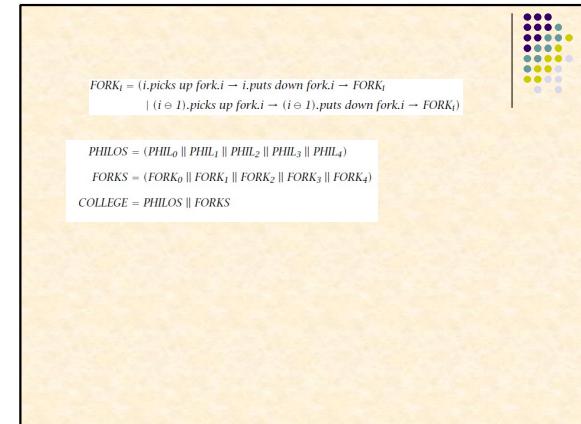
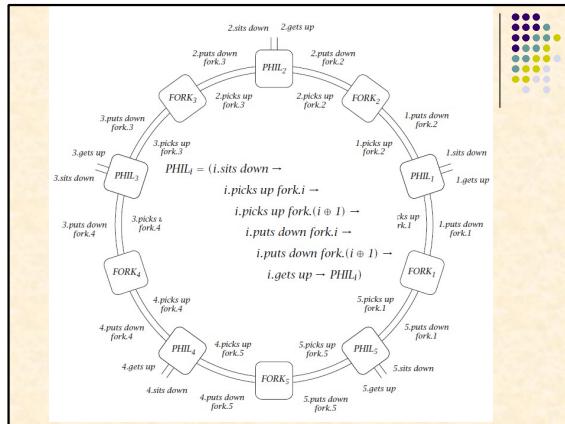
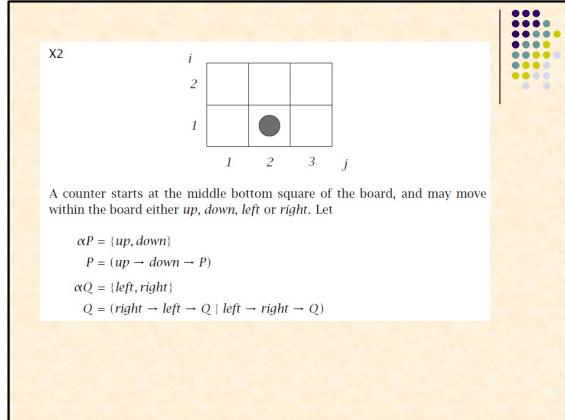
L4 $(x : A \rightarrow P(x)) \parallel (y : B \rightarrow Q(y)) = (z : (A \cap B) \rightarrow (P(z) \parallel Q(z)))$

Example

X1 Let $P = (a \rightarrow b \rightarrow P \mid b \rightarrow P)$
and $Q = (a \rightarrow (b \rightarrow Q \mid c \rightarrow Q))$
Then

$$\begin{aligned} (P \parallel Q) &= \\ &= a \rightarrow ((b \rightarrow P) \parallel (b \rightarrow Q \mid c \rightarrow Q)) \quad [\text{by L4A}] \\ &= a \rightarrow (b \rightarrow (P \parallel Q)) \quad [\text{by L4A}] \\ &= \mu X \bullet (a \rightarrow b \rightarrow X) \quad [\text{since the recursion is guarded.}] \end{aligned}$$

□



3.2.1 Laws

The algebraic laws governing nondeterministic choice are exceptionally simple and obvious. A choice between P and P is vacuous

L1 $P \sqcap P = P$ (idempotence)

It does not matter in which order the choice is presented

L2 $P \sqcap Q = Q \sqcap P$ (symmetry)

A choice between three alternatives can be split into two successive binary choices. It does not matter in which way this is done

L3 $P \sqcap (Q \sqcap R) = (P \sqcap Q) \sqcap R$ (associativity)

The occasion on which a nondeterministic choice is made is not significant. A process which first does x and then makes a choice is indistinguishable from one which first makes the choice and then does x

L4 $x \rightarrow (P \sqcap Q) = (x \rightarrow P) \sqcap (x \rightarrow Q)$ (distribution)

L5 $(x : B \rightarrow (P(x) \sqcap Q(x))) = (x : B \rightarrow P(x)) \sqcap (x : B \rightarrow Q(x))$

L6 $P \parallel (Q \sqcap R) = (P \parallel Q) \sqcap (P \parallel R)$

L7 $(P \sqcap Q) \parallel R = (P \parallel R) \sqcap (Q \parallel R)$

L8 $f(P \sqcap Q) = f(P) \sqcap f(Q)$

However, the recursion operator is *not* distributive, except in the trivial case where the operands of \sqcap are identical. This point is simply illustrated by the difference between the two processes

$P = \mu X \bullet ((a \rightarrow X) \sqcap (b \rightarrow X))$

$Q = (\mu X \bullet (a \rightarrow X)) \sqcap (\mu X \bullet (b \rightarrow X))$

4.2 Input and output

Let v be any member of $\alpha c(P)$. A process which first outputs v on the channel c and then behaves like P is defined

$(c!v \rightarrow P) = (c.v \rightarrow P)$

The only event in which this process is initially prepared to engage is the communication event $c.v$.

A process which is initially prepared to input any value x communicable on the channel c , and then behave like $P(x)$, is defined

$(c?x \rightarrow P(x)) = (y : \{ y \mid \text{channel}(y) = c \} \rightarrow P(\text{message}(y)))$

Example

X1 Using the new definitions of input and output we can rewrite 1.1.3 X7

$\text{COPYBIT} = \mu X \bullet (\text{in?}x \rightarrow (\text{out?}x \rightarrow X))$

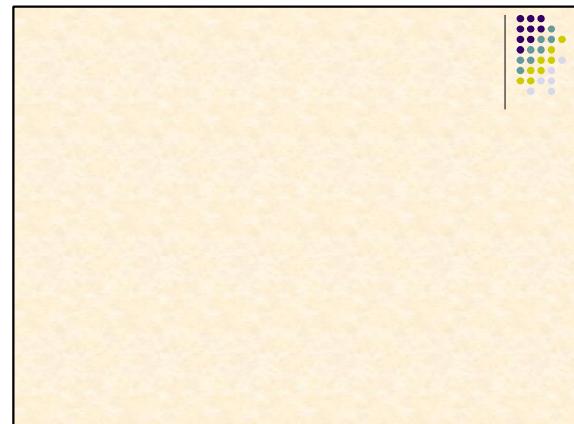
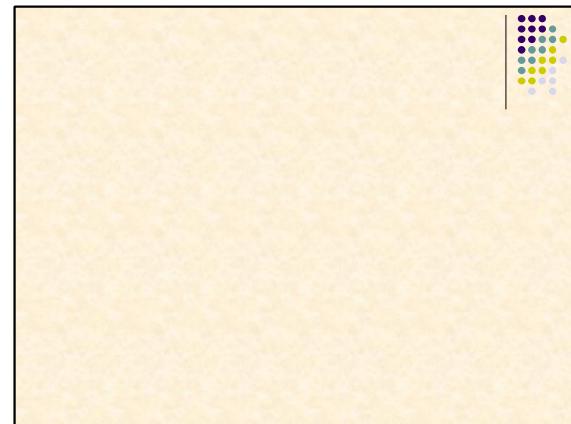
where $\text{ain}(\text{COPYBIT}) = \text{aout}(\text{COPYBIT}) = \{0, 1\}$ \square

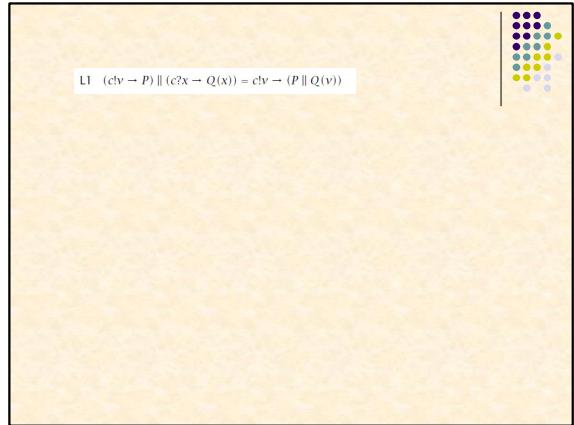
We shall observe the convention that channels are used for communication in only one direction and between only two processes. A channel which is used only for output by a process will be called an output channel of that process; and one used only for input will be called an input channel. In both cases, we shall say loosely that the channel name is a member of the alphabet of the process.

When drawing a connection diagram (Section 2.4) of a process, the channels are drawn as arrows in the appropriate direction, and labelled with the name of the channel (Figure 4.1).

Figure 4.1

L2 $((c!v \rightarrow P) \parallel (c?v \rightarrow Q(x))) \setminus C = (P \parallel Q(v)) \setminus C$
where $C = \{ c.v \mid v \in \alpha c \}$





TEMPORAL LOGIC

Heiko Krumm

University of Dortmund, Department of Computer Science

Symbolic logic generally supports the reasoning with propositions, i.e., with statements to be evaluated to true or false. Temporal logic is a special branch of symbolic logic focussing on propositions whose truth values depend on time. That contrasts with the classical logic point of view where the truth value of a repeatedly uttered proposition must always be the same and must neither depend on the modalities of the repetition nor on additional information. Temporal propositions typically contain some (explicit or implicit) reference to time conditions, while classical logic deals with timeless propositions. For instance consider the following examples:

A: “The moon is a satellite of the earth”

B: “The moon is rising”

C: “The moon is setting”

Proposition A can be viewed as timeless, since it is true in past, present, and future. In contrast, the propositions B and C have a temporalized aspect and refer to the implicit time condition “now”. Consequently temporal logic applies to time-related universes of discourse where behaviors and courses of events are of interest. As classical logic formulas can characterize static states and properties, temporal logic formulas can describe sequences of state changes and properties of behaviors.

Classical logic comprises different logics; several variants of propositional logic, first order predicate logic, etc., exist with different sets of logical operators and inference rules. Likewise some temporal logics were proposed which differ with respect to their formula syntax, semantics, and expressiveness. A temporal logic, however, basically results from an extension of a classical propositional or predicate logic by temporal quantifiers introducing temporalized modalities. Usually, there are at least the two quantifiers \blacksquare (denoting “always”) and \blacklozenge (denoting “eventually”) and typical formulas are similar to following examples:

D: $\blacklozenge B$ “The moon will be rising eventually”

E: $\blacksquare \blacklozenge B$ “The moon will be rising again and again”

F: $\blacksquare (B \Rightarrow \blacklozenge C)$ “Moon rise leads to moon setting”

The example formula D is true, if the moon is rising now or will be rising in some future point of time. Formula E exemplifies that combinations of temporal quantifiers can denote more complex time conditions, e.g., “always eventually” can correspond to the natural language term “again and again”. Finally, formula F is an example of a “leads-to” pattern describing that always a precondition B will eventually result in a postcondition C.

Due to its temporal quantifiers temporal logic is a convenient and appropriate means to reason with time-related propositions. Indeed, classical logic can also handle temporal properties, but the formulas tend to be complicated since points of time have to be explicitly represented in the underlying universe. The formula E may serve as example and underpin the usefulness of temporal logics. The easy-to-read temporal logic formula E corresponds to following predicate logic formula: “For all subjects x a subject y exists such, that – if x is a point of time – y is a point of time equal or later to x and the moon is rising at y ”.

History

Temporal logic is rooted in the field of exact philosophy and is a variant of modal logic. Modal logic deals with propositions which are interpreted with respect to a set of possible worlds. The truth value of propositions depends on the respective world and basically two operators “necessarily” and “possibly” exist which denote that a proposition is true in all possible worlds res. in some possible worlds. Even the ancient Greek philosophy schools of the Megarians, Stoics, and Peripatetics as well as Aristotle used some temporalized form of these modal operators. During the Middle Ages Arabian and European logicians resumed and refined the ancient approaches in order to discern different types of necessity and possibility. In modern times, the interest in symbolic logic grew during the first half of the 20th century, and – with some delay – new modal and temporal logic approaches occurred. First publications date back to the 1940's.

In particular, the logicians Prior, Rescher, Kripke, and Scott contributed to the development of modern temporal logic. Kripke presented a formal possible world semantics for modal logics. Prior proposed a temporal interpretation. An ordered set of possible worlds can correspond to a temporal sequence of states. In result, the two basic modal operators “necessarily” and “possibly” become the temporal quantifiers “always” and “eventually”. Based on the linearity of time additional operators like “next” and “until” as well as past operators were introduced. Rescher and Urquhart outlined the history and introduced several major systems of temporal logic in [5]. In 1974, Burstall proposed the application of

temporal logic in computer science for the first time. Pnueli improved this approach in [4], which is regarded as the classic source of temporal logic based program specification and verification.

Computer Science Application

In several fields of computer science there is a needs for the formal description of event-discrete processes and the corresponding reasoning. In the main, we have to mention the formal specification and verification of so-called reactive systems, the formalization of real-life processes as well as the semantics of natural language commands to be modeled in artificial intelligence, and finally the handling of dynamic consistency conditions in data base systems.

We focus on reactive systems. In particular, Manna and Pnueli recognized in [3] that reactive systems are of growing interest and that temporal logic is well-suited for their formal specification and verification. In contrast to those programs which transform starting states into final results and which may be specified by pre- and postconditions, reactive systems interact with their environment during runtime and the course of interactions and system events is essential. The range of reactive systems is wide and growing. It comprises embedded systems, process control systems, and all types of interactive, concurrent or distributed hard- and software systems. Due to their inherent concurrency, their elaborated fault-tolerance, coordination, and interaction mechanisms distributed systems are rather complex reactive systems and usually need particular design and development tools which support the formal handling of dynamic aspects. Here, temporal logic is profitably applied with respect to following topics:

1. Formal specification: Temporal logic formulas serve as precise, concise and binding descriptions of systems and components (e.g., as proposed by Lamport, Manna, and Pnueli in [2] res. [3]).
2. Formal verification: The rules of a temporal logic proof calculus are applied to show the correctness of a temporal logic specification with respect to more abstract system specifications (e.g., in [2] and [3]).
3. Requirements description: During the early system design the results of the requirements constraining the functional system behavior are represented by a set of temporal logic formulas.
4. Specification checks: Even if the design specifications use other means than temporal logic (e.g., other formal description techniques, *see* SDL, Estelle, and LOTOS, *see also* UNITY), temporal logic may be applied additionally in order to describe requirements and plausibility conditions. Meanwhile several approaches exist which support the tool-based checking of formal system specifications with respect to temporal logic conditions (*see* Model Checking).

Linear and Branching Time

Usually, a temporal logic can be classified as so-called linear-time logic which considers behaviors modeled as linear sequences of states. Within one behavior, each state has exactly one future. Additionally, so-called branching-time logics are known. Here, the formulas refer to tree-structured behaviors where a state can have several futures. The behavior-trees can directly correspond to tree models of non-deterministic systems (e.g., synchronization and communication trees, *see* Calculus of Communicating Systems). A corresponding prominent branching-time logic is CTL (computation tree logic, proposed by Clarke, Emerson, and Sistla in [1]). Its temporal quantifiers directly support the navigation in behavior trees.

Non-deterministic systems, however, have not necessarily to be modeled by behavior trees. Likewise, a set of linear state sequences can form a model of a non-deterministic system where each state sequence corresponds to one possible evolution of the system. In comparison with this linear-time approach, branching-time logics additionally provide for notions of potential behaviors since branching-time formulas can describe properties of branches which correspond to subsets of the possible execution sequences while linear-time formulas generally state properties of all possible sequences.

Variants

Besides of the mentioned distinctions between temporal propositional and predicate logics and between linear-time and branching-time logics, there exist further variants. Some introduce additional temporal quantifiers like “always in the past”, “sometimes in the past”, “next”, “precedes”, “until”, and “leads-to”. Others extend the time model, e.g., in order to describe time-intervals or real-time quantifications. Furthermore, partial-order temporal logics were proposed which directly refer to partial-order representations of concurrency (*see* Concurrency Model).

Example

To exemplify the application of temporal logic for the specification and verification of systems we outline some formula and proof patterns proposed by Lamport in [2] with respect to the Temporal Logic of Actions TLA which is a compact linear-time logic for the reasoning on state-transition systems. He considers the two commonly known classes of properties, invariance and eventuality. Moreover, the correctness of design refinements can be proven with respect to the preservation of properties.

An invariance property P is expressed by a formula “ $\blacksquare P$ ” where P is a predicate logic formula describing a set of execution states. Inter alia P may specify following typical correctness conditions of a system:

1. Partial correctness: P is an implication of the form “system terminated \Rightarrow correct results computed”.
2. Deadlock freedom: P applies to a set of states, the system is not deadlocked.
3. Mutual exclusion: P asserts that at most one process is in a critical section.

By means of auxiliary history variables all interesting safety properties of a system can be expressed as invariance properties (see Safety Property).

The formal proof of invariance properties is supported by a proof rule applying induction on the course of system execution steps. At first, one proves that each initial state implies P . Furthermore, each transition class of the system has to be considered. Each transition has to transform states where P is true into successor states where P is true again.

Eventuality properties assert that some events will eventually happen during each execution of a system. The following typical properties can be easily expressed in temporal logic:

1. Termination: A formula of the form “ \lozenge terminated” can assert that each execution leads to a state where the system is terminated.
2. Live service: Each state representing that a service request is pending will be followed by a state the request is served: “ $\blacksquare (\text{requested} \Rightarrow \lozenge \text{served})$ ”.
3. Fair message transfer: If a message is sent often enough over a loose channel, then it is eventually delivered: “ $(\blacksquare \lozenge \text{sent}) \Rightarrow (\lozenge \text{delivered})$ ”.

Eventuality properties can express the typical liveness requirements of systems (see Liveness Property).

The proof of eventuality properties can be reduced to the proof of a series of transitive leads-to properties of the form “ $\blacksquare (P \Rightarrow \lozenge Q)$ ”. The proof of a single leads-to property is supported by the so-called lattice rule which is based on the existence of a well-founded order. The order asserts that a finite number of execution steps is sufficient to reach a state where Q is true.

Systems can be described by formulas on abstract levels as well as on more implementation-near ones. Thus, specifications, refinement steps of a design, and implementations can be represented. That is of great interest, since valid implications correspond to system refinements which are correct in the usual understanding of system developers. Let the formula $Spec$ describe a system S on a more abstract level. A formula $Impl$ describes a correct refinement of S , if the implication formula “ $Impl \Rightarrow Spec$ ” is provable.

References

- [1] E.M. Clarke, E.A. Emerson, and A.P. Sistla, *Automatic Verification of Finite State Concurrent Systems Using Temporal Logic Specifications*, ACM Transactions on Programming Languages and Systems, 8(2): 244-263, 1986
- [2] L. Lamport, *The Temporal Logic of Actions*, ACM Transactions on Programming Languages and Systems, 16(3):872-923, 1994
- [3] Z. Manna and A. Pnueli, *The Temporal Logic of Reactive and Concurrent Systems*, Springer-Verlag, 1992
- [4] A. Pnueli, *The Temporal Logic of Programs*, Proceedings of the 18th IEEE Symposium on Foundations of Computer Science, pp. 46-57, 1977
- [5] N. Rescher and A. Urquhart, *Temporal Logic*, Springer-Verlag, 1971

Cross Reference:

CTL *see* Temporal Logic

Formal Specification *see* Temporal Logic

Formal Verification *see* Temporal Logic

TLA *see* Temporal Logic

Dictionary Terms:

Concurrency Model

Model representing the global dynamics of a system which consists of concurrently acting components. Mainly, there are two types of concurrency models. Interleaving models induce a total temporal ordering of all component actions. Thus, the system is assumed to perform a global sequence of actions and the model reduces concurrency to non-determinism. In contrast, partial-order models represent the temporal independence of concurrent events directly. Concurrent events are not comparable with respect to the order of events.

Liveness Property

Property of a system concerning its dynamics and expressing that the system will eventually show a particular behavior within a finite period of time. Together with safety properties (*see* Safety Property) liveness properties can be used to characterize the principal functionality of distributed systems.

Safety Property

Property of a system concerning its dynamics and expressing that the system behavior never injures particular conditions, e.g., never enters forbidden states. Together with liveness properties (*see* Liveness Property) safety properties can be used to characterize the principal functionality of distributed systems.