# CHAPTER 21

# CUDA dynamic parallelism

**Statically assign conservative worst-case grid**

**Dynamically provide higher resolution where accuracy is required**
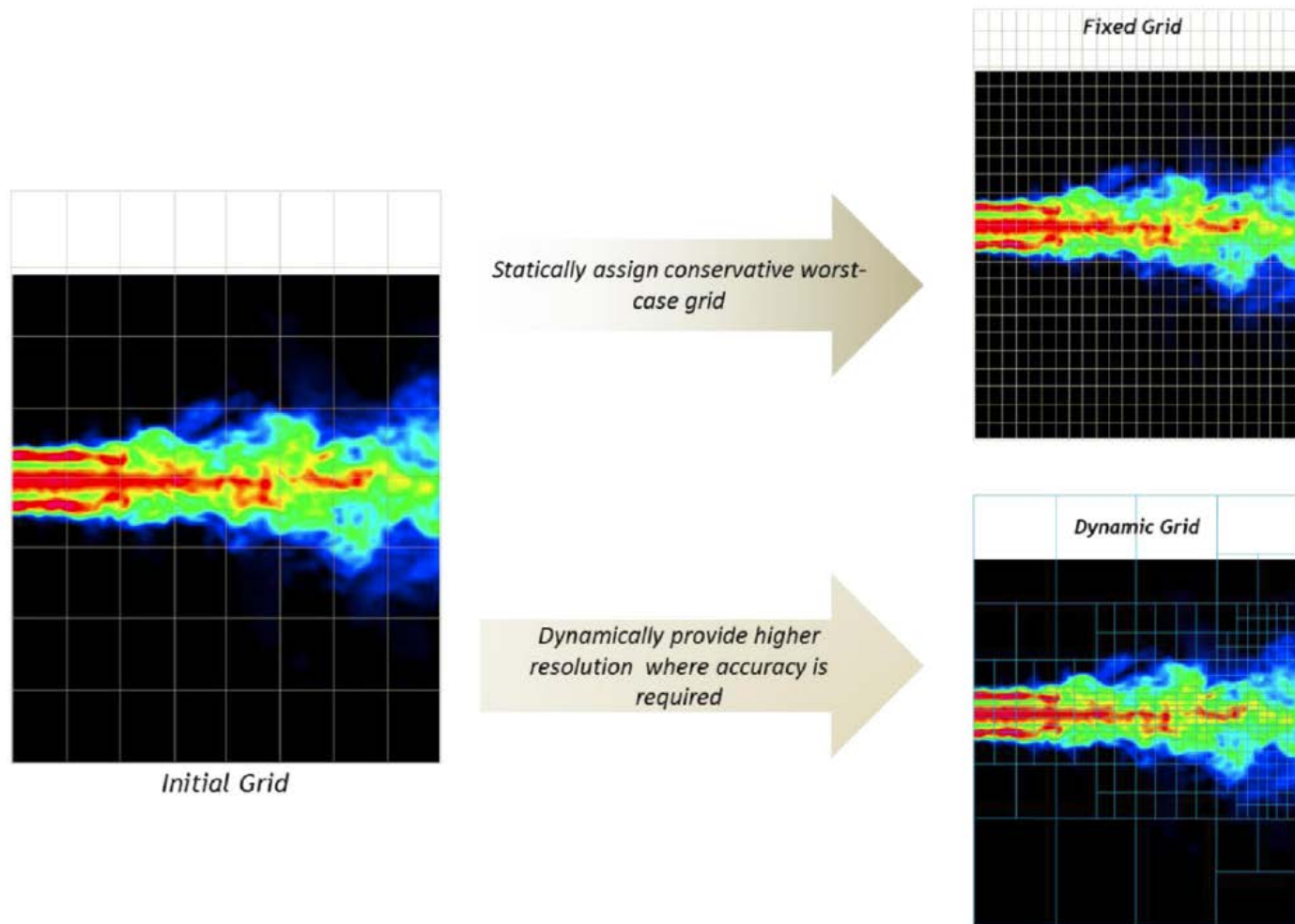
Fixed Grid

Dynamic Grid

Initial Grid

**FIGURE 21.1**
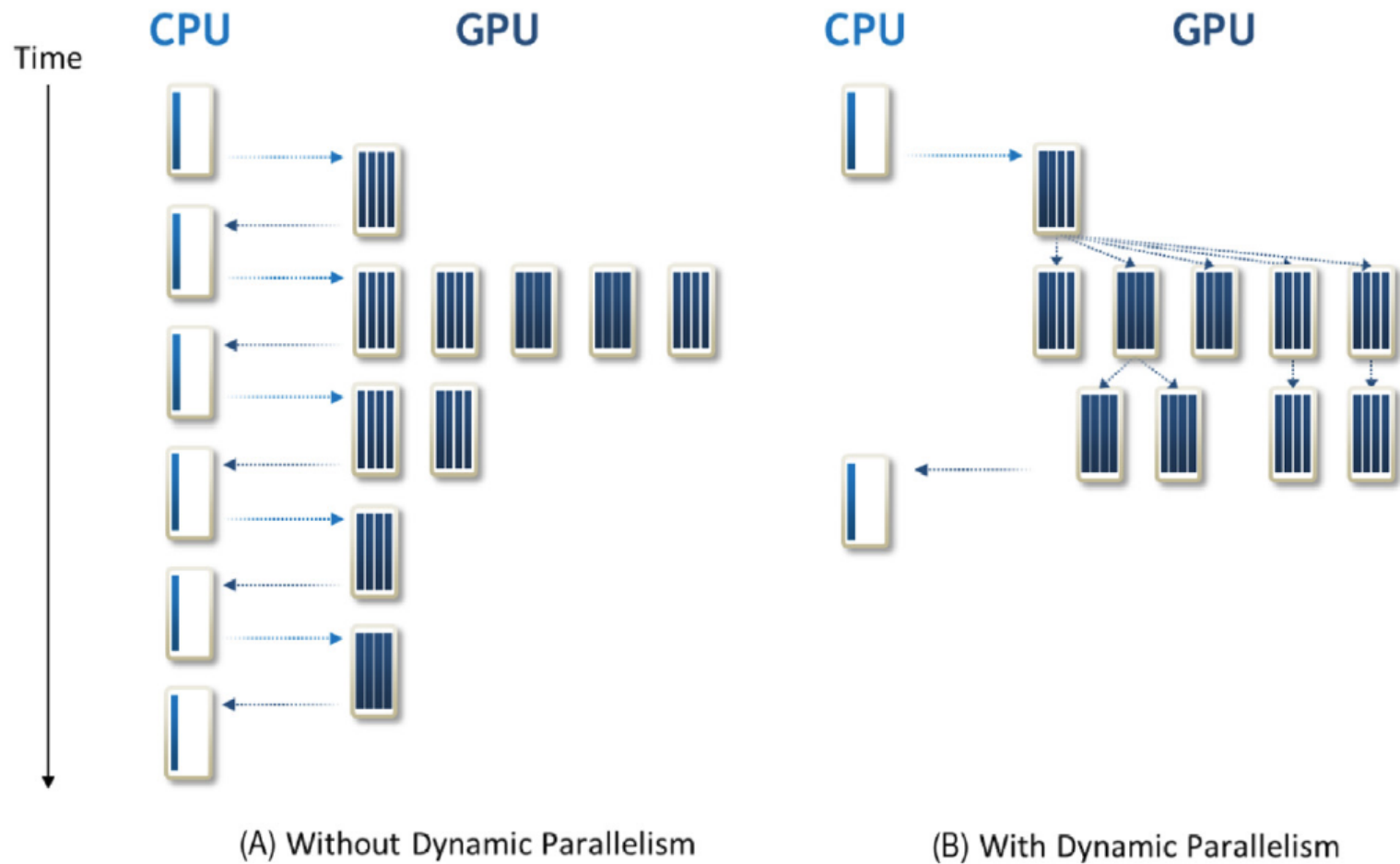
Fixed versus dynamic grids for a turbulence simulation model.

**FIGURE 21.2**

Grid launch patterns for algorithms with dynamic work variation, (A) without and (B) with dynamic parallelism.

```
int main() {
    float *data;
    setup(data);

    A <<< ... >>> (data);
    B <<< ... >>> (data);
    C <<< ... >>> (data);

    cudaDeviceSynchronize();
    return 0;
}
```

```
__global__ void B(float *data)
{
    do_stuff(data);

    X <<< ... >>> (data);
    Y <<< ... >>> (data);
    Z <<< ... >>> (data);
    cudaDeviceSynchronize();

    do_more_stuff(data);
}
```
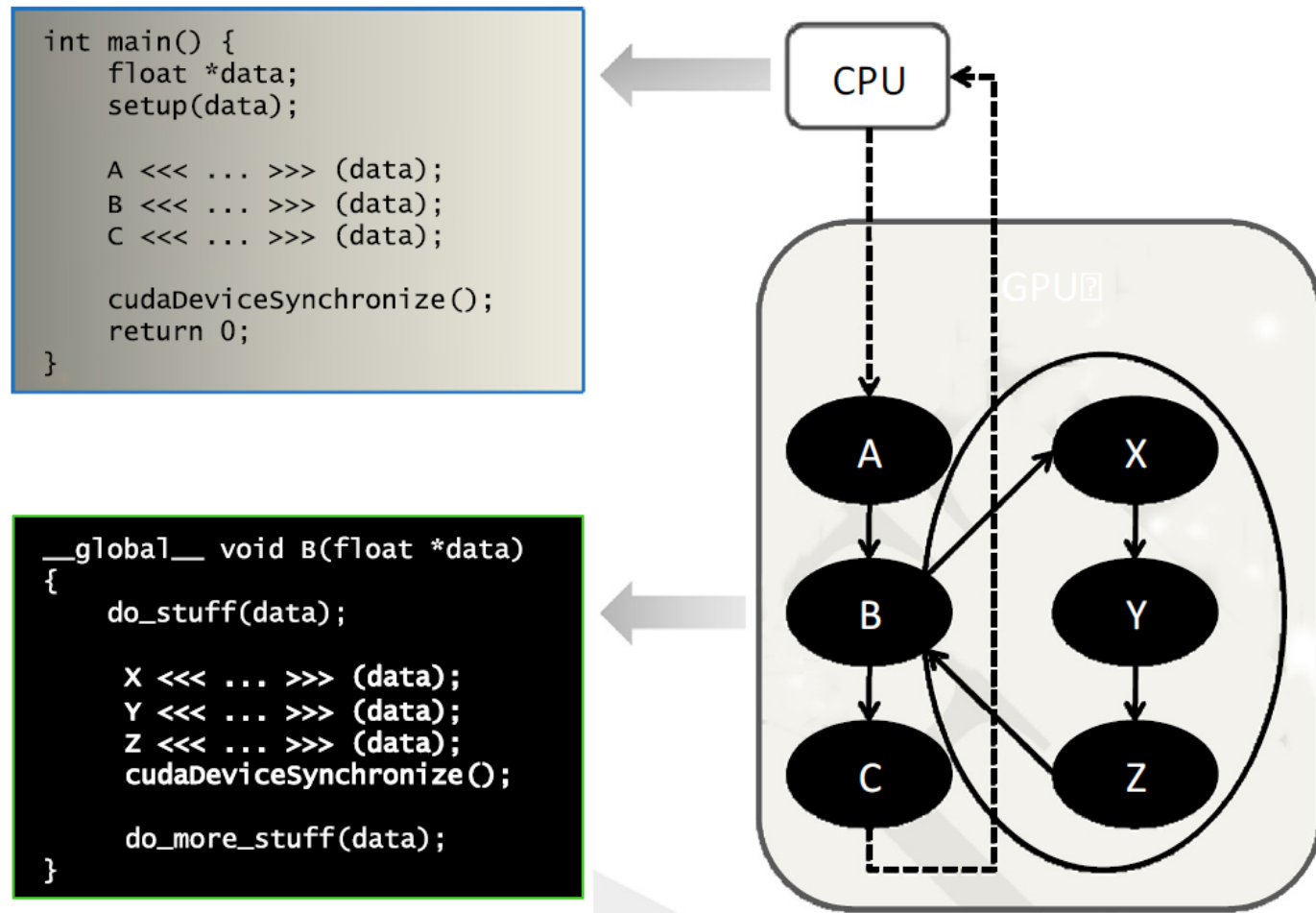
**FIGURE 21.3**

A simple example of a kernel (B) launching three kernels (X, Y, and Z).

```
01      __global__  void kernel(unsigned int* start, unsigned int* end,
02          float* someData, float* moreData) {
03
04          unsigned int i = blockIdx.x*blockDim.x + threadIdx.x;
05          doSomeWork_findMoreWork(someData[i]);
06
07          for(unsigned int j = start[i]; j < end[i]; ++j) {
08              doMoreWork(moreData[j]);
09          }
10
11      }
```

**FIGURE 21.4**

A simple example of a hypothetical parallel algorithm coded in CUDA without dynamic parallelism.

```
01      __global__ void kernel_parent(unsigned int* start, unsigned int* end,
02          float* someData, float* moreData) {
03
04          unsigned int i = blockIdx.x*blockDim.x + threadIdx.x;
05          doSomeWork(someData[i]);
06
07          kernel_child <<< ceil((end[i]-start[i])/256.0), 256 >>>
08              (start[i], end[i], moreData);
09
10      }
11
12      __global__ void kernel_child(unsigned int start, unsigned int end,
13          float* moreData) {
14
15          unsigned int j = start + blockIdx.x*blockDim.x + threadIdx.x;
16
17          if(j < end) {
18              doMoreWork(moreData[j]);
19          }
20
21      }
```

**FIGURE 21.5**

A revised example using CUDA dynamic parallelism.

```
01    #include <stdio.h>
02    #include <cuda.h>
03
04    #define MAX TESS POINTS 32
05
06    // A structure containing all parameters needed to tessellate a Bezier line
07    struct BezierLine {
08        float2 CP[3];                        //Control points for the line
09        float2 vertexPos[MAX_TESS_POINTS];   //Vertex position array to tessellate into
10        int nVertices;                       //Number of tessellated vertices
11    };
12
13    __global__ void computeBezierLines(BezierLine *bLines, int nLines) {
14        int bidx = blockIdx.x;
15        if(bidx < nLines){
16            //Compute the curvature of the line
17            float curvature = computeCurvature(bLines);
18
19            //From the curvature, compute the number of tessellation points
20            int nTessPoints = min(max((int)(curvature*16.0f),4),32);
21            bLines[bidx].nVertices = nTessPoints;
22
23            //Loop through vertices to be tessellated, incrementing by blockDim.x
24            for(int inc = 0; inc < nTessPoints; inc += blockDim.x){
25                int idx = inc + threadIdx.x;   //Compute a unique index for this point
26                if(idx < nTessPoints){
27                    float u = (float)idx/(float)(nTessPoints-1);   //Compute u from idx
28                    float omu = 1.0f - u;   //pre-compute one minus u
29                    float B3u[3]; //Compute quadratic Bezier coefficients
30                    B3u[0] = omu*omu;
31                    B3u[1] = 2.0f*u*omu;
32                    B3u[2] = u*u;
33                    float2 position = {0,0};   //Set position to zero
34                    for(int i = 0; i < 3; i++){
35                        //Add the contribution of the i'th control point to position
36                        position = position + B3u[i] * bLines[bidx].CP[i];
37                    }
38                    //Assign value of vertex position to the correct array element
39                    bLines[bidx].vertexPos[idx] = position;
40                }
41            }
42        }
43    }
```

**FIGURE 21.6**

Bezier curve calculation without dynamic parallelism.
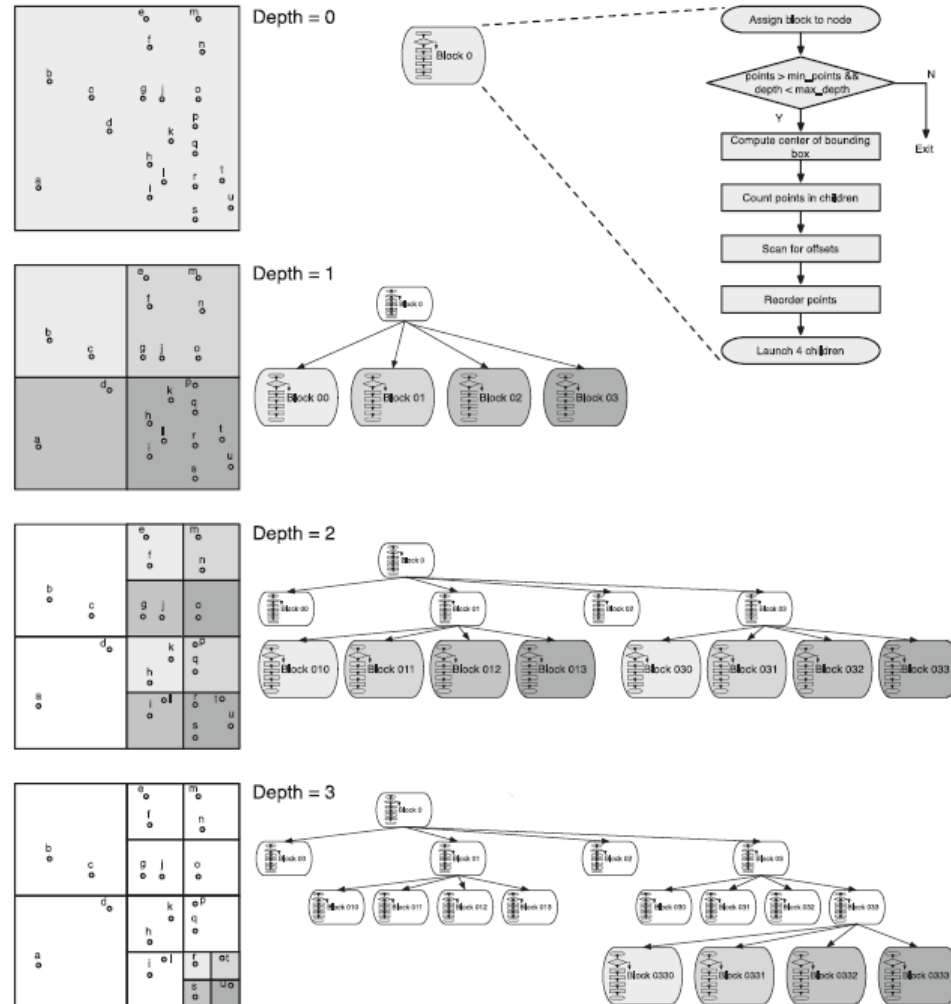
```
01    struct BezierLine {
02        float2 CP[3];       //Control points for the line
03        float2 *vertexPos; //Vertex position array to tessellate into
04        int nVertices;       //Number of tessellated vertices
05    };
06    __global__ void computeBezierLines_parent(BezierLine *bLines, int nLines) {
07        //Compute a unique index for each Bezier line
08        int lidx = threadIdx.x + blockDim.x*blockIdx.x;
09        if(lidx < nLines){
10            //Compute the curvature of the line
11            float curvature = computeCurvature(bLines);
12
13            //From the curvature, compute the number of tessellation points
14            bLines[lidx].nVertices = min(max((int)(curvature*16.0f),4),MAX TESS POINTS);
15            cudaMalloc((void**)&bLines[lidx].vertexPos,
16                bLines[lidx].nVertices*sizeof(float2));
17
18            //Call the child kernel to compute the tessellated points for each line
19            computeBezierLine_child<<<ceil((float)bLines[lidx].nVertices/32.0f), 32>>>
20                (lidx, bLines, bLines[lidx].nVertices);
21        }
22    }
23    __global__ void computeBezierLine child(int lidx, BezierLine* bLines,
24        int nTessPoints) {
25        int idx = threadIdx.x + blockDim.x*blockIdx.x;//Compute idx unique to this vertex
26        if(idx < nTessPoints){
27            float u = (float)idx/(float)(nTessPoints-1);  //Compute u from idx
28            float omu = 1.0f - u;    //Pre-compute one minus u
29            float B3u[3];    //Compute quadratic Bezier coefficients
30            B3u[0] = omu*omu;
31            B3u[1] = 2.0f*u*omu;
32            B3u[2] = u*u;
33            float2 position = {0,0};  //Set position to zero
34            for(int i = 0; i < 3; i++) {
35                //Add the contribution of the i'th control point to position
36                position = position + B3u[i] * bLines[lidx].CP[i];
37            }
38            //Assign the value of the vertex position to the correct array element
39            bLines[lidx].vertexPos[idx] = position;
40        }
41    }
42    __global__ void freeVertexMem(BezierLine *bLines, int nLines) {
43        //Compute a unique index for each Bezier line
44        int lidx = threadIdx.x + blockDim.x*blockIdx.x;
45        if(lidx < nLines)
46            cudaFree(bLines[lidx].vertexPos);    //Free the vertex memory for this line
47    }
```

**FIGURE 21.7**

Bezier calculation with dynamic parallelism.

**FIGURE 21.8**

Quadtree example. Each block is assigned to one quadrant. If the number of points in a quadrant is more than two, the block launches four child blocks. Shadowed blocks are active blocks in each level of depth.
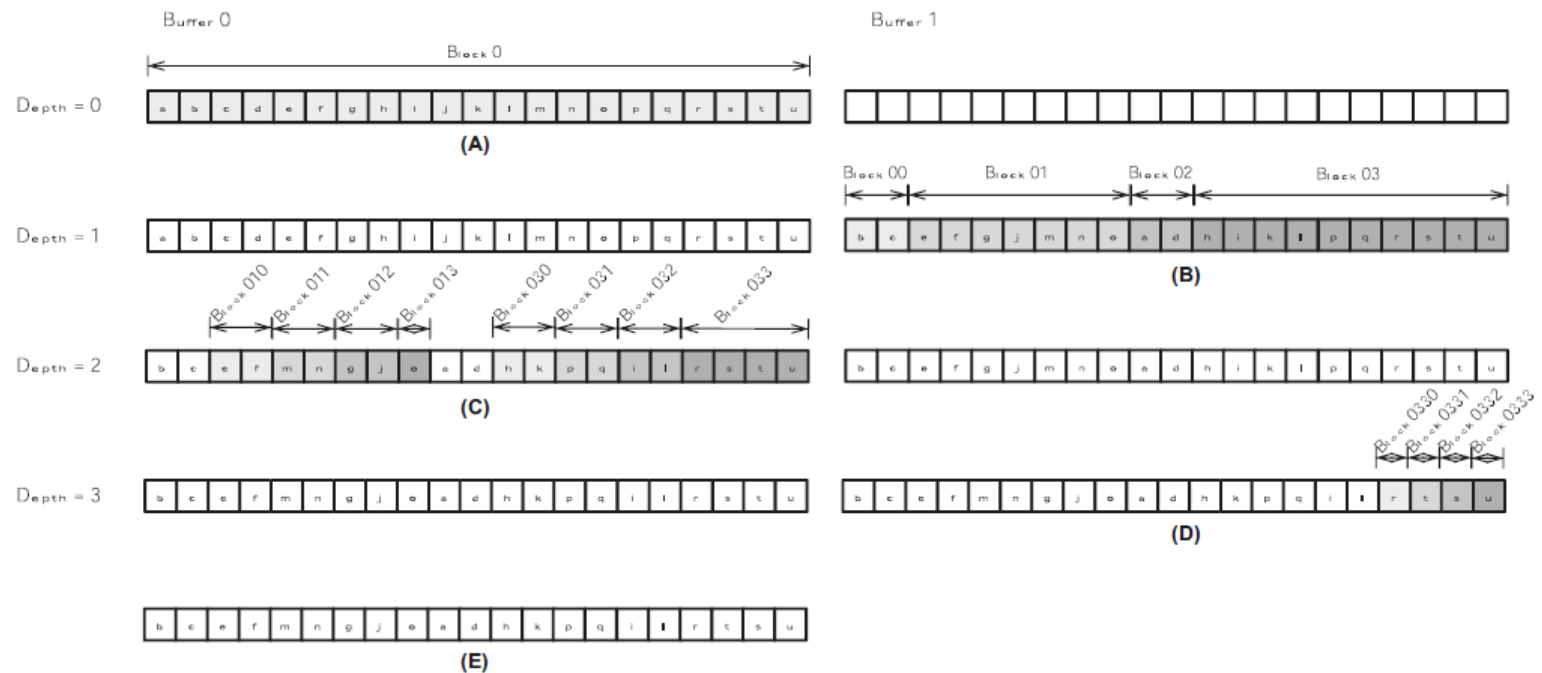
**FIGURE 21.9**

Quadtree example. At each level of depth, a block groups all points in the same quadrant together. (A) shows the initial input list in Buffer 0, (B) The list after being rearranged into four sublists that correspond to the four quadrants, (C) shows the list after being rearranged to reflect the second-level quadrants, (D) shows the list after being rearranged to reflect the third-level quadrant, (E) the final list is copied into Buffer 0 for return to the caller.

```
01    global    void build quadtree kernel
02             (Quadtree_node *nodes, Points *points, Parameters params) {
03      __shared__ int smem[8]; // To store the number of points in each quadrant
04
05      // The current node in the quadtree
06      Quadtree node &node = nodes[blockIdx.x];
07      node.set_id(node.id() + blockIdx.x);
08      int num_points = node.num_points(); // The number of points in the node
09
10      // Check the number of points and its depth
11      bool exit = check_num_points_and_depth(node, points, num_points, params);
12      if(exit) return;
13
14      // Compute the center of the bounding box of the points
15      const Bounding box &bbox = node.bounding_box();
16      float2 center;
17      bbox.compute_center(center);
18
19      // Range of points
20      int range_begin = node.points_begin();
21      int range_end   = node.points_end();
22      const Points &in_points = points[params.point_selector]; // Input points
23      Points &out_points = points[(params.point_selector+1) % 2]; // Output points
24
25      // Count the number of points in each child
26      count_points_in_children(in_points, smem, range_begin, range_end, center);
27
28      // Scan the quadrants' results to know the reordering offset
29      scan_for_offsets(node.points_begin(), smem);
30
31      // Move points
32      reorder_points(out_points, in_points, smem, range_begin, range_end, center);
33
34      // Launch new blocks
35      if (threadIdx.x == blockDim.x-1) {
36          // The children
37          Quadtree_node *children = &nodes[params.num_nodes_at_this_level];
38
39          // Prepare children launch
40          prepare_children(children, node, bbox, smem);
41
42          // Launch 4 children.
43          build quadtree kernel<<<4, blockDim.x, 8 *sizeof(int)>>>
44                      (children, points, Parameters(params, true));
45      }
46  }
```

**FIGURE 21.10**

Quadtree with dynamic parallelism: recursive kernel (support code in Appendix A210.1).

```
001   // Check the number of points and its depth
002     __device__  bool check_num_points_and_depth(Quadtree_node &node, Points *points,
003                                        int num_points, Parameters params){
004       if(params.depth >= params.max_depth || num_points <= params.min_points_per_node) {
005         // Stop the recursion here. Make sure points[0] contains all the points
006         if(params.point_selector == 1) {
007           int it = node.points_begin(), end = node.points_end();
008           for (it += threadIdx.x ; it < end ; it += blockDim.x)
009             if(it < end)
010               points[0].set_point(it, points[1].get_point(it));
011         }
012       return true;
013     }
014     return false;
015   }
016
017   // Count the number of points in each quadrant
018     __device__  void count_points_in_children(const Points &in_points, int* smem,
019       int range_begin, int range_end, float2 center) {
020       // Initizalize shared memory
021       if(threadIdx.x < 4) smem[threadIdx.x] = 0;
022         __syncthreads();
023       // Compute the number of points
024       for(int iter=range_begin+threadIdx.x; iter<range_end; iter+=blockDim.x){
025           float2 p = in_points.get_point(iter); // Load the coordinates of the point
026           if(p.x < center.x && p.y >= center.y)
027               atomicAdd(&smem[0], 1); // Top-left point?
028           if(p.x >= center.x && p.y >= center.y)
029               atomicAdd(&smem[1], 1); // Top-right point?
030           if(p.x < center.x && p.y < center.y)
031               atomicAdd(&smem[2], 1); // Bottom-left point?
032           if(p.x >= center.x && p.y < center.y)
033               atomicAdd(&smem[3], 1); // Bottom-right point?
034       }
035       __syncthreads();
036   }
037
038   // Scan quadrants' results to obtain reordering offset
039     __device__  void scan_for_offsets(int node_points_begin, int* smem){
040       int* smem2 = &smem[4];
041       if(threadIdx.x == 0){
042       for(int i = 0; i < 4; i++)
043           smem2[i] = i==0 ? 0 : smem2[i-1] + smem[i-1]; // Sequential scan
044       for(int i = 0; i < 4; i++)
045           smem2[i] += node_points_begin;  // Global offset
046       }
047         __syncthreads();
048   }
049
050   // Reorder points in order to group the points in each quadrant
051   __device__  void reorder_points(
```

**FIGURE 21.11**

Quadtree with dynamic parallelism: device functions (support code in Appendix A21.1).

```
052                      Points& out_points, const Points &in_points, int* smem,
053                      int range begin, int range end, float2 center){
054        int* smem2 = &smem[4];
055        // Reorder points
056        for(int iter=range begin+threadIdx.x; iter<range end; iter+=blockDim.x){
057            int dest;
058            float2 p = in_points.get_point(iter); // Load the coordinates of the point
059            if(p.x<center.x && p.y>=center.y)
060                dest=atomicAdd(&smem2[0],1); // Top-left point?
061            if(p.x>=center.x && p.y>=center.y)
062                dest=atomicAdd(&smem2[1],1); // Top-right point?
063            if(p.x<center.x && p.y<center.y)
064                dest=atomicAdd(&smem2[2],1); // Bottom-left point?
065            if(p.x>=center.x && p.y<center.y)
066                dest=atomicAdd(&smem2[3],1); // Bottom-right point?
067            // Move point
068            out points.set point(dest, p);
069        }
070        __syncthreads();
071    }
072
073    // Prepare children launch
074      device    void prepare children(Quadtree node *children, Quadtree node &node,
075                                       const Bounding box &bbox, int *smem){
076        int child_offset = 4*node.id(); // The offsets of the children at their level
077
078        // Set IDs
079        children[child offset+0].set id(4*node.id()+ 0);
080        children[child offset+1].set id(4*node.id()+ 4);
081        children[child offset+2].set id(4*node.id()+ 8);
082        children[child_offset+3].set_id(4*node.id()+12);
083
084        // Points of the bounding-box
085        const float2 &p min = bbox.get min();
086        const float2 &p max = bbox.get max();
087
088        // Set the bounding boxes of the children
089        children[child_offset+0].set_bounding_box(
090            p_min.x , center.y, center.x, p_max.y);      // Top-left
091        children[child offset+1].set bounding box(
092            center.x, center.y, p max.x , p max.y);      // Top-right
093        children[child offset+2].set bounding box(
094            p min.x , p min.y , center.x, center.y);      // Bottom-left
095        children[child_offset+3].set_bounding_box(
096            center.x, p_min.y , p_max.x , center.y);      // Bottom-right
097
098        // Set the ranges of the children.
099        children[child offset+0].set range(node.points begin(),   smem[4 + 0]);
100        children[child offset+1].set range(smem[4 + 0], smem[4 + 1]);
101        children[child_offset+2].set_range(smem[4 + 1], smem[4 + 2]);
102        children[child_offset+3].set_range(smem[4 + 2], smem[4 + 3]);
103    }
```

**FIGURE 21.11**

(Continued)

```
// Create non-blocking stream
cudaStream t stream;
cudaStreamCreateWithFlags(&stream, cudaStreamNonBlocking);

//Call the child kernel to compute the tessellated points for each line
computeBezierLine_child<<<ceil((float)bLines[lidx].nVertices/32.0f), 32, 0, stream>>>
    (lidx, bLines, bLines[lidx].nVertices);

// Destroy stream
cudaStreamDestroy(stream);
```

**FIGURE 21.12**

Child kernel launch with named streams.

```
kernel_name<<< Dg, Db, Ns, S >>>([kernel arguments])
```

In-text figure 1

```
__global__ void parent_kernel(int *output, int *input, int *size) {
    // Thread index
    int idx = threadIdx.x + blockDim.x*blockIdx.x;

    // Number of child blocks
    int numBlocks = size[idx] / blockDim.x;

    // Launch child
    child_kernel<<< numBlocks, blockDim.x >>>(output, input, size);

}
```

In-text figure 2

In-text figure 3

```
001    // A structure of 2D points
002    class Points {
003        float *m_x;
004        float *m_y;
005
006        public:
007        // Constructor
008          __host__    __device__   Points() : m_x(NULL), m_y(NULL) {}
009
010        // Constructor
011        __host__    __device__  Points(float *x, float *y) : m_x(x), m_y(y) {}
012
013        // Get a point
014        __host__    __device__   __forceinline__  float2 get_point(int idx) const {
015            return make_float2(m_x[idx], m_y[idx]);
016        }
017
018        // Set a point
019        __host__    __device__     __forceinline__   void set_point(int idx, const float2 &p) {
020            m_x[idx] = p.x;
021            m_y[idx] = p.y;
022        }
023
024        // Set the pointers
025        __host__    __device__   __forceinline__   void set(float *x, float *y) {
026            m_x = x;
027            m_y = y;
028        }
029    };
030
031    // A 2D bounding box
032    class Bounding_box {
033        // Extreme points of the bounding box
034        float2 m_p_min;
035        float2 m_p_max;
036
037        public:
038        // Constructor. Create a unit box
039        __host__    __device__   Bounding_box() {
040            m_p_min = make_float2(0.0f, 0.0f);
041            m_p_max = make_float2(1.0f, 1.0f);
042        }
043
044        // Compute the center of the bounding-box
045        __host__    __device__   void compute_center(float2 &center) const {
046            center.x = 0.5f * (m_p_min.x + m_p_max.x);
047            center.y = 0.5f * (m_p_min.y + m_p_max.y);
048        }
049
050        // The points of the box
051        __host__    __device__   __forceinline__   const float2 &get_max() const {
052            return m_p_max;
053        }
054
055        __host__    __device__     __forceinline__   const float2 &get_min() const {
056            return m_p_min;
057        }
058
059        // Does a box contain a point
060        __host__    __device__   bool contains(const float2 &p) const {
061            return p.x>=m_p_min.x && p.x<m_p_max.x && p.y>=m_p_min.y && p.y<m_p_max.y;
062        }
063
```

In-text figure 4

```
064        // Define the bounding box
065            __host__ __device__ void set(float min_x, float min_y, float max_x, float
max_y) {
066            m_p_min.x = min_x;
067            m_p_min.y = min_y;
068            m_p_max.x = max_x;
069            m_p_max.y = max_y;
070        }
071    };
072
073    // A node of a quadree
074    class Quadtree_node {
075        // The identifier of the node
076        int m_id;
077        // The bounding box of the tree
078        Bounding_box m_bounding_box;
079        // The range of points
080        int m_begin, m_end;
081
082    public:
083        // Constructor
084        __host__ __device__ Quadtree_node() : m_id(0), m_begin(0), m_end(0) {}
085
086        // The ID of a node at its level
087        __host__ __device__ int id() const {
088            return m_id;
089        }
090
091        // The ID of a node at its level
092        host    device   void set_id(int new_id) {
093            m_id = new_id;
094        }
095
096        // The bounding box
097        __host__ __device__ __forceinline__ const Bounding_box &bounding_box() const {
098            return m_bounding_box;
099        }
100
101        // Set the bounding box
102        host    device    forceinline   void set_bounding_box(float min_x,
103            float min_y, float max_x, float max_y) {
104            m_bounding_box.set(min_x, min_y, max_x, max_y);
105        }
106
107        // The number of points in the tree
108        __host__ __device__ __forceinline__ int num_points() const {
109            return m_end - m_begin;
110        }
111
112        // The range of points in the tree
113        host    device    forceinline   int points_begin() const {
114            return m_begin;
115        }
116
117        host    device    forceinline   int points_end() const {
118            return m_end;
119        }
120
121        // Define the range for that node
122        __host__ __device__ __forceinline__ void set_range(int begin, int end) {
123            m_begin = begin;
124            m_end = end;
125        }
126    };
127
128    // Algorithm parameters
129    struct Parameters {
```

```
130          // Choose the right set of points to use as in/out
131          int point selector;
132          // The number of nodes at a given level (2^k for level k)
133          int num_nodes_at_this_level;
134          // The recursion depth
135          int depth;
136          // The max value for depth
137          const int max depth;
138          // The minimum number of points in a node to stop recursion
139          const int min_points_per_node;
140
141          // Constructor set to default values.
142          __host__ __device__ Parameters(int max_depth, int min_points_per_node) :
143             point selector(0),
144             num_nodes_at_this_level(1),
145             depth(0),
146             max depth(max depth),
147             min_points_per_node(min_points_per_node) {}
148
149          // Copy constructor. Changes the values for next iteration
150          __host__ __device__ Parameters(const Parameters &params, bool) :
151             point_selector((params.point_selector+1) % 2),
152             num nodes at this level(4*params.num nodes at this level),
153             depth(params.depth+1),
154             max_depth(params.max_depth),
155             min points per node(params.min points per node) {}
156      };
```

In-text figure 5