

Production and Playback of Human Figure Motion for Visual Simulation

John P. Granieri, Jonathan Crabtree, Norman I. Badler

University of Pennsylvania, Philadelphia

Abstract

We describe a system for off-line production and real-time playback of motion for articulated human figures in 3D virtual environments. The key notions are (1) the logical storage of full-body motion in posture graphs, which provides a simple motion access method for playback, and (2) mapping the motions of higher DOF figures to lower DOF figures using slaving to provide human models at several levels of detail, both in geometry and articulation, for later playback. We present our system in the context of a simple problem: Animating human figures in a distributed simulation, using DIS protocols for communicating the human state information. We also discuss several related techniques for real-time animation of articulated figures in visual simulation.

CR Categories and Subject Descriptors: I.6.3 [Simulation and Modeling]: *Applications*; I.3.7 [Three-Dimensional Graphics and Realism]: *Animation*; E.1 [Data Structures]: *graphs*;

General Terms: Animation

Additional Key words and phrases: posture graphs, real-time animation, multi-resolution motion, visual simulation

1

¹ Authors' address: Center for Human Modeling and Simulation, University of Pennsylvania, Philadelphia, PA 19104-6389, USA

This work was presented in an earlier form during the Virtual Reality Annual International Symposium (VRAIS) '95, March 11-15 1995, sponsored by the IEEE.

This work was supported in part by ARO DAAL03-89-C-0031 including U.S. Army Research Laboratory; Naval Training Systems Center N61339-93-M-0843; Sandia Labs AG-6076; ARPA AASERT DAAH04-94-G-0362; DMSO DAAH04-94-G-0402; ARPA DAMD17-94-J-4486; U.S. Air Force DEPTH through Hughes Missile Systems F33615-91-C-0001; DMSO through the University of Iowa; and NSF CISE CDA88-22719.

1 Introduction

The ability to render realistic motion is an essential part of many virtual environment and visual simulation applications. Nowhere is this more true than in virtual worlds containing simulated humans. Whether these human figures represent the users' virtual personae (avatars) or computer-controlled characters, people's innate sensitivity as to what looks "natural" with respect to human motion demands, at the very least, that moving characters be updated with each new frame that the image generator produces.

We first discuss a topical problem in visual simulation which requires the real-time rendering of realistic human motion, and then describe our system for authoring the motion off-line, and playing back that motion in real time. We also address some of the issues in real-time image generation of highly-articulated figures, as well as look at several other methods used for real-time animation.

2 Human motion in DIS

The problem we are interested in is generating and displaying motion for human figures, in particular soldiers, in a distributed virtual environment. Parts of the general problem and the need for representing simulated soldiers (referred to as Dismounted Infantry, or DIs), are covered in [21, 5]. Although primarily driven by military requirements today, the general technologies for projecting real humans into, and representing simulated humans within, virtual environments, should be widely applicable in industry, entertainment and commerce in the near future.

The Distributed Interactive Simulation (DIS) [9] protocol is used for defining and communicating human state information in the distributed virtual environment.

A typical distributed simulation contains many simulation hosts, each concerned with simulating a portion, or sub-set, of all the objects (or *entities*) involved in a simulation (here, entity can refer to a human figure, a vehicle, or other part of the environment) and processes involved in the simulation.

DIS defines a protocol for heterogeneous simulation applications to interoperate, typically for the real-time simulation of battlefield operations. It defines the packets of information (which are referred to as protocol data units - PDUs) and the set of rules for exchanging the packets between simulation applications, with the goal of achieving a shared and correlated synthetic environment between the applications. For a good example of the software structure inside a simulation application in this environment, the reader is referred to [13].

A simple example configuration of a distributed simulation is shown in Fig. 1. Simulator A runs a simulation of a soldier entity moving across a terrain (we are not particularly interested in how that is done in this paper, but an example system would be the SAFDI system [17]). As the soldier moves in simulator

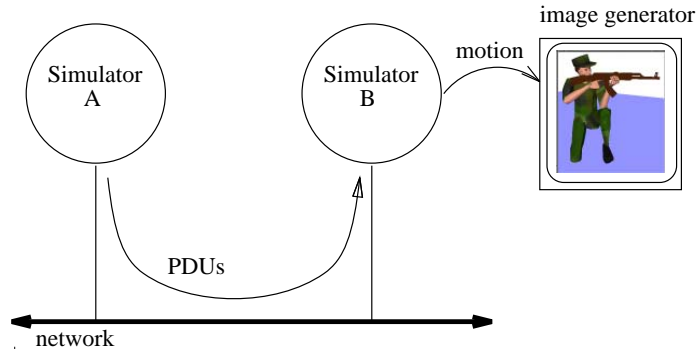


Figure 1: A simple distributed simulation

A, his essential state information is propagated over time to simulator B via a stream of Entity State PDUs. Simulator B monitors this stream and updates its internal state information concerning the soldier. Simulator B also has an image generator as a component, to visually display the simulated environment and entities. In this paper, we are interested in that part of the system in simulator B which is used to create animations of the soldier figure in the image generator, according to the stream of PDUs coming from simulator A. Note that simulator B may be simulating a set of its own soldiers, whose state information would be propagated back to simulator A, via the same mechanisms. The techniques used in this paper for animating a human figure in a DIS-based distributed simulation introduce several issues relating to latency and coherence between the soldier's state information in simulator A and B. The interested reader is referred to [16].

The DIS protocol, at least the part relating to human entities, is in its early stages of development, and fairly limited in what it can specify about a human figure [16], but is a good baseline to start with.

The information representing a human entity is currently defined by several discrete enumerations in the appearance field of an Entity State Protocol Data Unit (PDU) in the DIS protocol [10]. The relevant information we are interested in from the Entity State PDU is shown in Fig. 2. The human is always in one of the four postures, along with a weapon state. The heading defines the forward direction. Although there are enumerations for walking and crawling, we use combinations, such as (*posture=standing*)+(velocity>0) to be equivalent to walking or running. Although the protocol allows for up to three weapons of different types on a soldier, we only modeled one, a rifle.

If the human can be in any of n possible postures, there are potentially n^2 transitions between the postures. Rather than create n^2 posture transitions, we encode the postures and transitions into a *posture graph* [1]. The graph defines the motion path to traverse to move the human figure from any one posture to

| Field | Value | Units |
|----------|---|---------------|
| Posture | <i>Standing</i> <i>Kneeling</i> <i>Prone</i> <i>Dead</i> | n/a |
| Weapon | <i>Deployed</i> <i>Firing</i> | n/a |
| Position | P_x, P_y, P_z | meters |
| Velocity | V_x, V_y, V_z | meters/second |
| Heading | <i>theta</i> | degrees |

Figure 2: Essential human state information in a DIS Entity State PDU

another. These graphs are directed and may include cycles. It also provides the logical structure for the run-time motion database.

The posture graph we are using is shown in Fig. 3, partitioned into two sub-graphs for facilitating discussion (the dotted lines around the Standing Deployed and Prone Deployed nodes merely represent duplicate nodes in the diagram – they are the same as the similarly labeled nodes with solid lines, attached via the dotted line).

When the velocity of the human is zero, the possible transitions between *fixed* postures are encoded in the right-hand posture graph of Fig. 3. We use the term “fixed” as it is descriptive of when this graph is used: typically the soldier figure transitions to one of the fixed postures, represented in a node, and stays fixed in that posture for some time.

A few of the actual fixed postures are shown in Fig. 4.

When the velocity of the figure is non-zero, the possible transitions between *locomotion* postures are shown in the left-hand posture graph of Fig. 3. In this graph, the nodes are static postures, but the figure would never be in the posture for more than one frame.

In the posture graph, the nodes represent postures, and the directed arcs represent the animated full-body transitions, or movements, from posture to posture. Each arc stores the intermediate postures, or motion, and has an associated time for traversal. The time associated with the transition arc is used to find the shortest path, in time, if more than one path exists between a starting posture and a goal posture. The human figure is always in a posture defined in a node, or in one of the intermediate postures defined along the arcs. A useful analogy for visualizing the graph is shown in Fig. 5. Three strips of film are pasted together, representing three arcs of the graph converging on a node (the shared frame in the middle). In the left sequence, the stick figure is raising its left arm (its looking at you). In the upper sequence, it is raising its

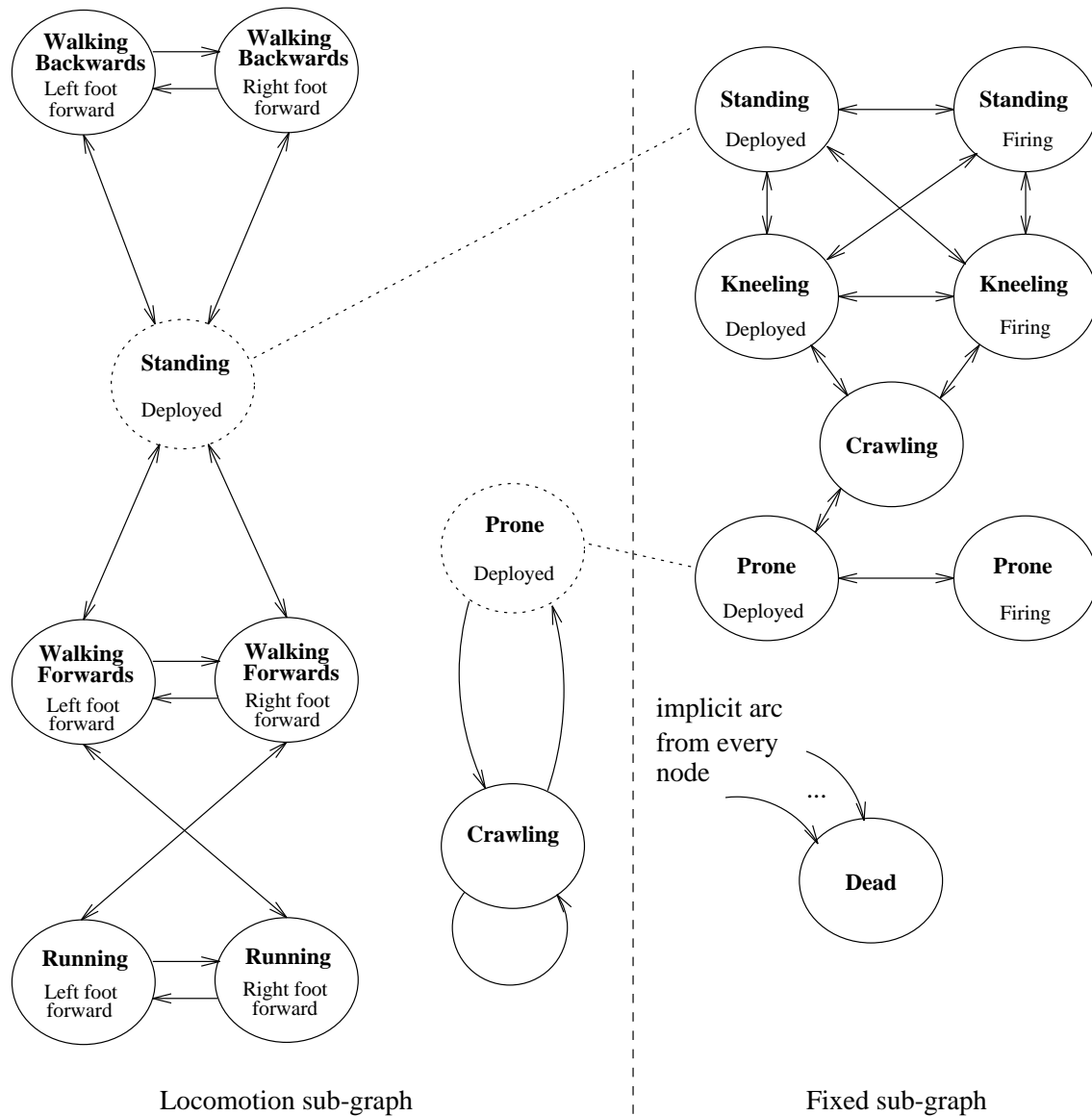


Figure 3: The posture graph. The dotted lines represent duplicated nodes, to facilitate discussion. The dashed line partitions the graph into the left sub-graph, referred to as the *locomotion* graph, and the right sub-graph, referred to as the *fixed* posture graph

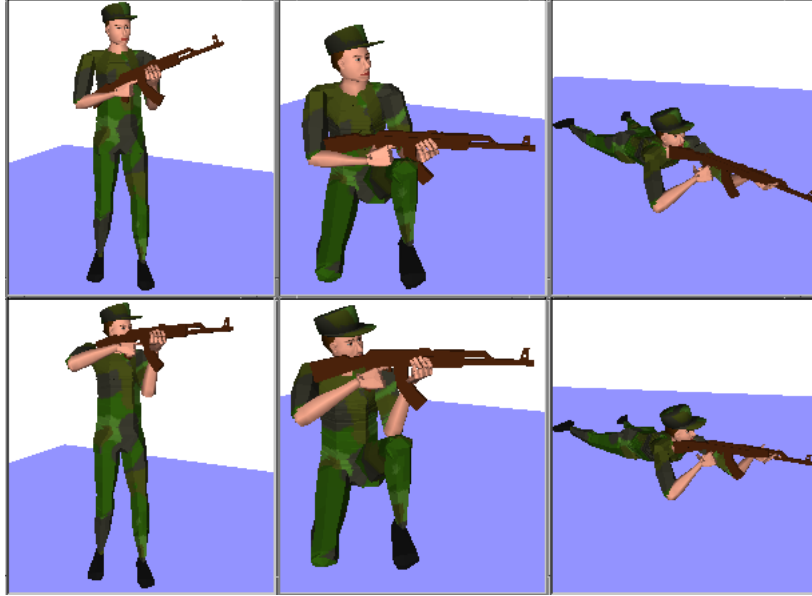


Figure 4: Some of the postures a soldier can take in DIS

right foot, and in the right sequence, it is waving its left arm. The motion on individual arcs of the graph is similar to the notion of animation *tracks* described in [6]. Here we are putting many of the tracks together, with common starting and ending frames, to form the posture graph.

The system we built consists of two distinct parts: 1) the off-line motion data generator, and 2) the on-line real-time playback mechanism, running in a high-performance IRIS Performer-based [18] image generator application.

3 Off-line motion production

Motion production involves three steps: 1) creating postures and motion for each node and arc in a posture graph, for one human model, 2) mapping the resulting motion onto human models with lower degrees-of-freedom (DOF) and lower resolution geometry, and finally 3) recording the results and storing in a format for easy retrieval during playback.

3.1 Authoring the motion

The first step in producing motion for real-time playback is to create postures representing the nodes in the posture graphs, as well as the corresponding motions between them, represented as the directed arcs in the graphs. We used a

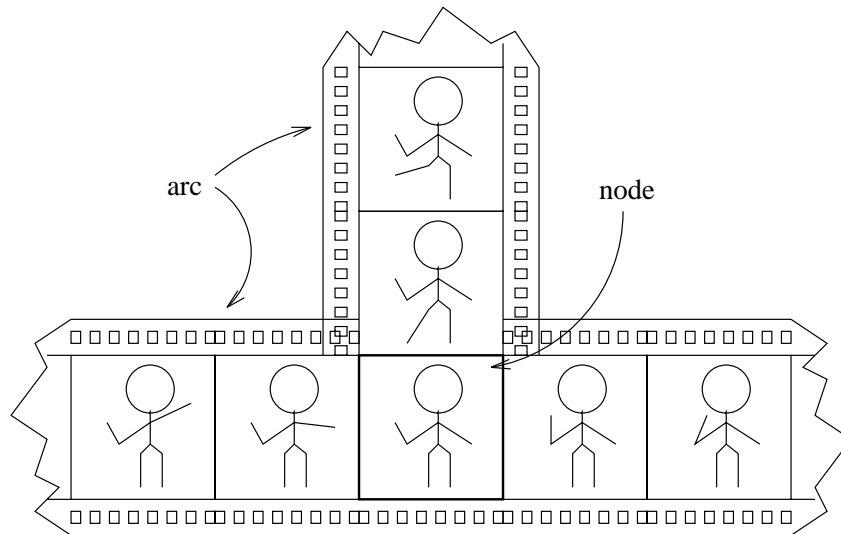


Figure 5: Nodes and arcs of a posture graph seen as strips of film

slightly modified version of the *Jack* human modeling and animation system [2] for this purpose. *Jack* provides a nice constraint-based, goal-driven system (relying heavily on inverse-kinematics and primitive “behavioral” controls) for animating human figures, as well as facilities for organizing motions for general posture interpolation [1]. It is important to note that the posture graphs presented in this paper differ from the *posture transition graphs* presented in [1]. In the latter, the posture transition graphs are used to organize motion primitives for general posture interpolation with collision avoidance. In the former application (this paper) the posture graphs are a logical mechanism for organizing a database of pre-recorded motion, and determining motion sequences as paths between nodes of the graph. An underlying assumption of the posture graphs is that the articulated human figure’s motion is continuous, and therefore can be organized into a connected graph.

Each directed transition in the *fixed* posture graph typically was produced from 10 to 15 motion primitives (e.g. `move_arm`, `bend_torso`, etc). Many of the directed motions from a posture node A to a posture node B are simply run in reverse to get the corresponding motion from posture B to posture A. In several cases, the reverse motion was scripted explicitly for more natural results.

The human figure can also move (either forwards or backwards, depending on the difference between the heading and the direction of the velocity vector) by either locomoting (if posture is standing) or crawling (if posture is prone). The locomotion posture graph transitions of Fig. 3 were generated by Hyeongseok Ko’s walking system [11]. Six strides for each type of walking transition were

| | human-high | human-med | human-low |
|----------------|------------|-----------|-----------|
| polygons | 2400 | 500 | 120 |
| rigid segments | 69 | 19 | 12 |
| joints | 73 | 17 | 11 |
| DOFs | 134 | 50 | 21 |
| motion | 60Hz | 30Hz | 15Hz |

Figure 6: The different levels of detail for the human models

generated (forward walking, backward walking, running): left and right starting steps, left and right ending steps, and left and right cyclic steps. The crawling animation was generated manually, and is based on two animations - one that goes from the prone posture to the cyclic state, and one complete cyclic motion. Note that only straight line locomotion of fixed stride is modeled. We are currently working on extending the system to handle variable stride length and curved path locomotion, as possible in the system of [11].

3.2 Slaving

The second step in the production process is concerned with preparing the motion for the real-time playback system. We wish to have tens, and potentially hundreds, of simulated humans in a virtual environment. This necessitates having multiple level-of-detail (LOD) models, where the higher resolution models can be rendered when close to the viewpoint, and lower resolution models can be used when farther away (see Section 4.1 for the reasoning). We reduce the level of detail in the geometry and articulation by creating lower-resolution (both in geometry and articulation) human figures, with the characteristics listed in the table of Fig. 5.

All the motions and postures of the first step are authored on a (relatively) high resolution human body model which includes fully articulated hands and spine. This model is referred to as “human-high” in the above table. We manually created the two lower-resolution models, human-med and human-low. Because of the difference in internal joint structure between human-high and the lower LOD models, their motions cannot be controlled by the available human control routines in *Jack* (which all make assumptions about the structure of the human figure - they assume a structure similar to human-high). Instead of controlling their motion directly, we use the motion scripts generated in the first step to control the motion of a human-high, and then map the motion onto the lower resolution human-med and human-low. We call this process *slaving*, because the high resolution figure acts as the *master*, and the low resolution figure acts as the *slave*.

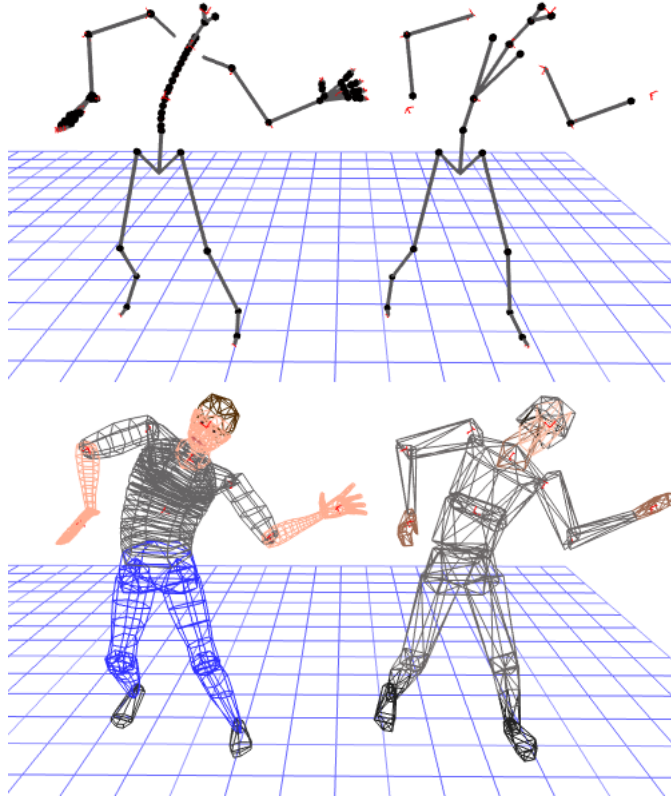


Figure 7: human-high and human-med models during slaving. human-high is the master. Upper window is the skeletal articulation. Models are offset for illustrative purposes.

Even though the different LOD human models have different internal joint structures and segment shapes, their gross dimensions (e.g., length of arms, torso, etc.) are similar. The slaving process consists of interpolating the motions for the full human figure, generating all the in-between frames, and simultaneously having a lower LOD human model (human-med or human-low) slaved, and then saving the in-between frames for the soldier. We will describe the process used for slaving from human-high to human-med; the case with human-low is similar.

For each frame of an animation, we first compute the position and joint angles for human-high. Then, an approximation of the joint angles for human-med are computed. This is straightforward, as certain joints are the same (the elbow, for example, is only one DOF on both human models), and others can

be approximated by linear combinations (for example, the 35 DOFs of the spine on human-high can be summed and mapped directly onto the 7 DOF torso of human-med). This gives a good first approximation of the posture mapping, and provides an initial configuration for the final mapping. For the resulting motion of human-med to look correct, we need to have certain landmark sites of the two bodies match exactly (the hands must be on the rifle). The final mapping step involves solving a set of constraints (point-to-point and orientation), to bring the key landmark sites into alignment. The constraints are solved using an iterative inverse kinematics routine [23] to move the body parts into alignment.

Because of differences in geometry between the master and slave, in general we cannot expect all the landmark sites to match exactly. For the problem domain of this paper, animating a soldier figure from a stream of DIS Entity State PDUs, the hands are always holding a rifle, so matching the hand positions accurately from the master is very important (otherwise the slave's hands may penetrate the rifle). Using a priority scheme in evaluating constraints, we assign higher priority to the hand-matching constraints than others, to account for this fact. If the slaving procedure cannot fit the master and slave within a certain tolerance, it will generate a warning for the animator.

3.3 Recording

The final step in the motion production process is to record the resulting motions of the human figures. The recorded motion for one transition is referred to as a *channel set* (where each joint or figure position is referred to as a *channel*; the channel is indexed by time). Again, the channel set is analogous to an animation track in [6]. For each LOD human figure, a homogeneous transform is recorded, representing figure position relative to a fixed point, and for each joint, the joint angles are recorded (one angle per DOF). Also for joints, the composite joint transform is pre-computed and stored as a 4x4 matrix (which can be plugged directly into the parenting hierarchy of the visual database of the run-time system). Each channel set has an associated transition time. The channels of human-high are interpolated and stored at 60Hz, human-med at 30Hz, and human-low at 15Hz. These rates correspond to the motion sampling during playback (see below).

4 Real-time motion playback

The real-time playback functions are packaged as a single linkable library, intended to be embedded in a host IRIS Performer-based visual simulation application. The library loads the posture graphs shown in Fig. 4, as well as the associated channel set motion files. Only one set of motions are loaded, and shared amongst any number of soldier figures being managed by the library. The articulated soldier figures themselves are loaded into the Performer run-time vi-

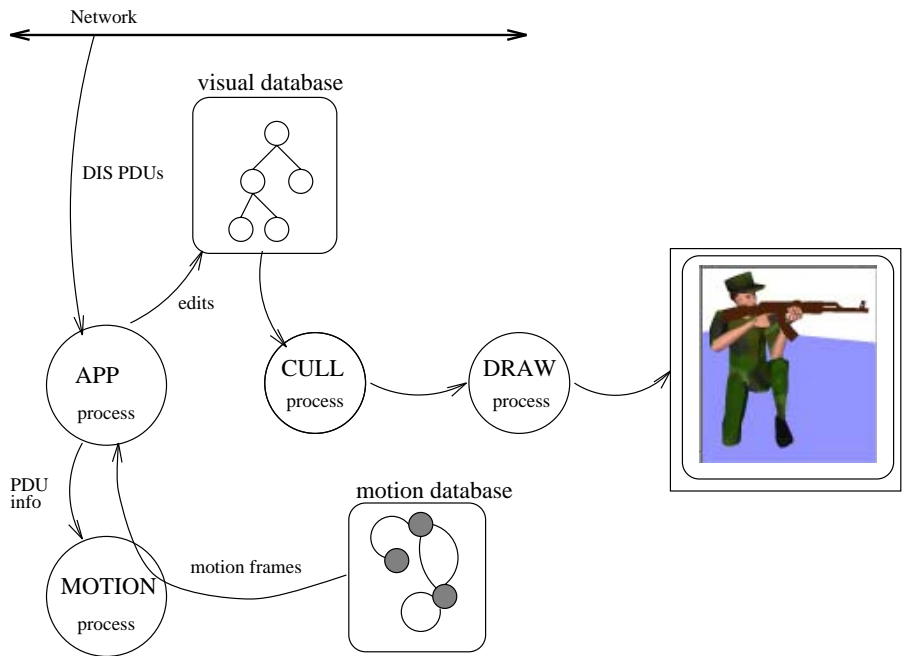


Figure 8: Overview of multi-processing framework for run-time system. Circles are processes, rounded-squares are data structures, and arrows are data flows

sual database. The library runs as a separate process, the MOTION process, serving motion data to the APP process (the APP, CULL and DRAW process are defined in the Performer multiprocessing framework, and form a software rendering pipeline, which increases rendering throughput, and also latency). See Fig. 8 for a schematic overview of the runtime system. The APP process is the main controlling process, and would be where any simulation processing would occur. For a good description of what goes on in the APP process of a typical DIS-based simulation application, the reader is referred to [13].

An update function in the APP process is provided which maps joint angle values into the joint transforms of the soldier figures in the Performer visual database.

The APP process receives PDUs containing state information for remotely-simulated soldier figures. It extracts the essential information (as shown in Fig. 2) and sends this as requests to the MOTION process. The MOTION process translates these requests into playbacks of channel sets (the traversal of arcs of the posture graphs), and continuously sends joint angle packets back to the APP process. The APP process then performs the necessary edits to the articulation matrices in the visual database, and then initiates a frame to be rendered. The CULL and DRAW processes form the software rendering pipeline, where the CULL

culls the visual database to a viewing frustum, and then passes this modified database to the DRAW process, which efficiently feeds it to the head of the hardware graphics rendering pipeline, which then renders the final image to the frame buffer.

In the case of a fixed posture change (a motion between any two nodes in the fixed posture graph of Figure 3), the system will find the shortest path (as defined by traversal time) between the current and goal postures in the graph, and execute the sequence of transitions. For example, if the posture graph is currently at Standing Deployed, and Prone Firing is requested, it will transition from Stand Deployed to Crawl to Prone Deployed, and finally to Prone Firing. The Dead posture node in Fig. 3 is special in that there is an implicit arc from every node in the graph to this node. The motion along these implicit arcs is not recorded and stored, as along all other arcs, but generated on the fly via simple joint interpolations.

A posture change is made with a node of the fixed graph as a destination only upon receipt of a DIS Entity State PDU indicating that the agent is in such a posture. In the absence of further information, the agent **remains** in that posture. Conversely, when a posture change is made with a node of the locomotion graph as the destination, something that will occur if a PDU indicates the agent now has a non-zero speed, the agent does **not** remain in that posture once it is reached; absence of further information in this case means that the agent's speed is still nonzero, and hence the agent must take another step, or crawl another meter forwards, or whatever is appropriate for the current mode of locomotion. This continued motion requires that another posture change be made immediately.

A simple finite-state controller manages the posture transitions within the locomotion graph. For instance, the transition from Standing Deployed to Walking Forwards (left foot forward) is taken whenever the agent's speed becomes non-zero and the agent's heading vector agrees with the velocity vector. On the other hand, if the vectors are not pointing in approximately the same direction, a transition is instead made to one of the Walking Backwards states. While the agent's speed remains nonzero (as it is assumed to in the absence of PDU updates), the controller continually makes transitions back and forth between, for example, the Walking Forwards (left foot forward) and Walking Forwards (right foot forward) nodes. This cycle of transitions creates a smooth walking motion by concatenating successive left and right steps. Note that since we currently have no cycles of more than two nodes, finding the shortest path between postures in a cycle is a trivial matter! Crawling is handled similarly, though it is a simpler case; there is no need for separate "left foot forward" and "right foot forward" states.

The system samples all the pre-recorded motion using elapsed time, so it is guaranteed to always play back in real time. For a 2 second posture transition recorded at 60fps, and a current frame rate of the image generator of 20fps, the playback system would play frames 0, 3, 6, ..., 120. It recomputes the elapsed

transition time on every frame, in case the frame rate of the image generator is not uniform.

The motion frame update packets sent from the MOTION process back to the APP process are packaged to include only those joint angles which have changed from the last update. This is one way we can minimize joint angle updates, and take advantage of frame-to-frame coherence in the stored motions². A full update (all joint angles and figure positions) is about 400 bytes.

4.1 Motion level-of-detail

It is recognized that maintaining a constant frame rate is essential to the believability of a simulation, even if it means accepting an update speed bounded by the most complex scene to be rendered. Automatic geometric level-of-detail selection, such as that supported by the IRIS Performer toolkit, is a well-known technique for dynamically responding to graphics load by selecting the version of a model most appropriate to the current viewing context [4, 7, 20].

The LOD selection within the visual database seeks to minimize polygon flow to the rendering pipeline (both in the software CULL and DRAW components of the software pipeline, as well as to the transformation engines within the hardware pipeline).

Given our representation, which enforces the separation of geometry and motion, it is possible to expand level of detail selection into the temporal domain, through *motion level-of-detail* selection. In addition to reducing polygon flow, via selecting lower LOD geometric models, we are also selecting lower LOD articulation models, with fewer articulation matrices, as well as sampling motion at lower frequencies. This reduces the flow of motion updates (joint angles and matrices) to the articulation matrices in the visual database. The models we are using are listed in Fig. 3.2. For example, a motion sequence lasting 1 second, rendered on human-high, would generate $(73 \text{ joints})(1 \text{ second})(60 \text{ updates/sec}) = 4380 \text{ jointupdates}$. On human-low, the same motion sequence generates $(11 \text{ joints})(1 \text{ second})(15 \text{ updates/sec}) = 165 \text{ jointupdates}$.

Techniques which employ data represented at various resolutions are sometimes referred to as *multi-resolution* or *multi-scale* techniques or analyses. When working with images, the resolution or scale factor is usually two (i.e. an image would be represented at 256x256, 128x128, 64x64, etc), and when working with successive levels-of-detail in geometry, the scale factor is an order-of-magnitude (i.e. a geometric model at 1000 triangles, 100 triangles, and 10 triangles). We chose a scale factor of 2 (60Hz, 30Hz, 15Hz) between our recorded motions. There wasn't a tremendous amount of deep thought behind the choice, but it does seem to work well for the current application. A more rigorous selection

²An initial implementation of the playback library was run as an independent process, on another machine, from the host image generator, and joint angle packets were sent over TCP/IP stream sockets, hence the desire to minimize net traffic.

of frequencies, taking in to account rules of sampling, should be done in the future.

In the playback system, we simultaneously transition to a different geometric representation with a simpler articulation structure, and switch between stored motions for each articulation model. We gain performance in the image generator, while consuming more run-time storage space for the motions. Our metric for LOD selection is simply the distance to the virtual camera. This appears to work satisfactorily for our current application domain, but further evaluation of the technique, as well as more sophisticated selection metrics (e.g. the metrics described in [7, 4]) need to be explored.

5 Example implementations and uses

The real-time playback system has been used in two DIS-based applications to create motion for simulated soldiers in virtual environments.

The Team Tactical Engagement Simulator [21] projects one or more soldiers into a virtual environment, where they may engage hostile forces and practice coordinated team activities. See Fig. 9 for a sample view into the training environment. The soldier stands in front of a large projection screen, which is his view into the environment. He has a sensor on his head and one on his weapon. He locomotes through the environment by stepping on a resistive pad and controls direction of movement and field of gaze by turning his head. The soldier may also move off the movement pad, and the view frustum is updated accordingly based on his eye position (head-coupled display). This allows the soldier, for example, to crouch down to see under a parked vehicle, or to peek around the corner of a building while still affording himself the protection of the building. TTES also creates the necessary DIS Entity State PDUs to represent the real soldier (mapping from sensor values into the small set of postures in the Entity State PDU), and sends them out over the net to other TTES stations that are participating in the exercise. Currently, TTES uses an extension to the system as described in this paper. Initially, they used authored motions for the posture graphs. They currently use motion-capture data, from work done by Boston Dynamics Inc. [15] (joint angle motion from real soldiers moving through the postures, captured via attached magnetic or optical sensors, and then cleaned up with post-processing filters).

The playback system is also used in a version of the NPSNET-IV [5] system, for generating motion of SIMNET- and DIS-controlled soldier entities, as well as the motions for battlefield medics.

Motion level-of-detail selection is of particular relevance to the example projects described above, because in the situation where a hostile agent enters the field of view of a soldier (one of the real human participants) and brings his weapon into the deployed position, the hostile's actions will probably be noted only in the participant's peripheral vision. It is well-known that humans



Figure 9: A View of Battle Town with several soldiers in different postures

can detect the presence of motion in their peripheral vision very easily, but that resolution of detail is very low. When head-tracking data is available or a head-mounted display is in use, it is possible to designate areas of the viewing frustum as peripheral and reduce geometric and motion detail accordingly (not just based on linear distance to the camera, but angular offsets also). In the TTES environment this “focus of attention” information can be obtained from the aim of the real soldier’s rifle when it is in the raised position, as the real soldier will almost certainly be sighting in this situation.

The system has also been integrated with a behavioral simulation which can navigate a pedestrian agent about a simple urban environment, while respecting sidewalks, crosswalks, streetlights, and other pedestrians. The behavioral simulation is decoupled from the motion generation; it simply schedules steps for the pedestrian, and the motion system described here creates a smoothly locomoting human figure [8].

6 Comparison of production/playback methods

One of the most obvious criteria for evaluating a given motion representation is size; there is a clear progression in the methods used to animate humans (or any entity whose geometric representation varies over time) based on the amount of space required to store a given motion. We look at three methods.

The first method, requiring the most storage, involves generating and rendering the movements of characters in an off-line fashion. Frame-by-frame, a sequence of two-dimensional snapshots is captured and saved for later playback. The image generator then displays the bit-mapped frames in sequence,

possibly as texture maps on simple rectangular polygons. Hardware support for texture mapping and alpha blending (for transparent background areas in the texture bitmaps) make this an attractive and fast playback scheme. Furthermore, mip-mapping takes care of level-of-detail management that must be programmed explicitly in other representations. Since the stored images are two-dimensional, it is frequently the case that artists will draw each frame by hand. In fact, this is precisely the approach utilized in most video games for many years. It is clear that very little computation is required at run-time, and that altering the motions incurs a high cost and cannot be done in real time. In fact, modifying almost any parameter except playback speed must be done off-line, and even playback speed adjustments are limited by the recording frequency. However, one real problem with using two-dimensional recording for playback in a three-dimensional scene is that non-symmetric characters will require the generation of several or many sets of frames, one for each possible viewing angle, increasing storage requirements still further. The authors of the popular game DOOM [19] record eight views of each animated character (for each frame) by digitizing pictures of movable models, and at run time the appropriate frames for the current viewing angle (approximately) are pasted onto a polygon. These eight views give a limited number of realistic viewing angles; it is impossible, for instance, to view a DOOM creature from directly above or below. Interestingly enough, an article on plans for a follow-up to DOOM reveals that the authors intend to switch to one of the two remaining representations we describe here:

Unlike the previous games, the graphic representation of characters will be polygon models with very coarse texture mapping. This will make it hard to emulate natural locomotion, so they'll stay away from creating too many biped characters.[22]

Making the move to the second method involves a relatively slight conceptual change, namely taking 3-dimensional snapshots instead of 2-dimensional snapshots. This means storing each frame of a figure's motion as a full three-dimensional model. Doing so obviates the need for multiple data sets corresponding to multiple viewing positions and shifts slightly more of the computational burden over to the image generator. Instead of drawing pixels on a polygon the run-time system sends three-dimensional polygonal information to the graphics subsystem. It is still an inflexible approach because the figures are stored as solid "lumps" of geometry (albeit textured), from which it is extremely difficult, if not impossible, to extract the articulated parts of which the original model is comprised. Modifications must still be effected off-line, although rendering is done in real time. This is essentially the approach used by the SIMNET image generators to display soldiers on a simulated battlefield [3].

The final method is the one implemented by the system described in this paper, in which we record not the **results** of the motions, but the motions themselves. We store a single **articulated** three-dimensional model of each

agent, and from frame to frame record only the joint angles between articulated segments. Modern rendering toolkits such as the IRIS Performer system used in this project increasingly allow support for storing coordinate transformations within a visual database, with relatively little cost associated with updating the transformation matrices in real time. As a result of adopting this approach, storage space is reduced and it is far easier to accurately perform interpolation between key frames because articulation information is not lost during motion recording. It also allows for virtual agents with some motions replayed strictly “as-is” and some motions which may be modified or generated entirely in real time. For instance, the slight changes in shoulder and elbow joint orientation required to alter the aim of a weapon held by a virtual soldier could be generated on demand.

We believe that the smallest representation presented in our list, the third method, actually retains the **most** useful information and affords the most flexibility, while placing an acceptable amount of additional computational burden on the run-time display system.

7 Extensions & future work

We are currently exploring several extensions to the techniques described above, to add more expressive power to the tool bag of the real-time animator.

Key-framing and interpolation The use of the pre-recorded motions in the above posture graphs trades time for space. We do not compute joint angles on the fly, but merely sample stored motions. As the motions become more complex, it becomes very time-consuming to produce all the motions in the off-line phase, so we only produce key frames in a transition, every 5 to 10 frames, and then use simple interpolations to generate the inbetweens during real-time playback. This technique can’t be extended much beyond that, as full-body human motion does not interpolate well beyond that many frames. This also reduces the amount of stored motions by a factor proportional to the spacing of the key frames, but increases computation time when a playback frame lands between two key frames. Also, replacing the recorded motion in some transition arcs with purely procedural motion generators can further reduce storage. For example, a sub-graph of a posture graph which only contains a few joints moving (a left-handed wave) can easily be replaced with a simple function of time which returns the correct shoulder, elbow, and wrist angles during different phases of the transition.

Partitioning full-body motion In the posture graphs described previously, each motion transition included all the joint angles for the whole body. A technique to reduce motion storage, while increasing playback flexibility,

is to partition the body into several regions, and record motion independently for each region. For example, the lower body can be treated separately during locomotion, and the upper body can have a variety of different animations played on it. Also, to support the mapping of motion from partially sensed real humans (i.e. sensors on the hands) onto the animated human figures, we want to animate the lower body and torso separately, then place the hands and arms using a fast inverse kinematics solution.

Articulation level-of-detail The various LOD models we used for the human figures were all built manually. Techniques for synthesizing lower LOD geometric models are known, but they don't apply to building lower articulation LOD models. Some techniques for automatically synthesizing the lower articulation skeletal models, given a high resolution skeleton and a set of motions to render, would be very useful.

Other dynamic properties A limitation is currently imposed by the fact that the segments of our articulated figures must be rigid. However, this is more an implementation detail than a conceptual problem, since with sufficient computational power in the run-time system real-time segment deformation will become possible. In general it seems likely that the limiting factor in visual simulation systems will continue to be the speed at which the graphics subsystem can actually render geometry. The adoption of coarse-grained multiprocessing techniques [18] will allow such operations as rigid or elastic body deformations to be carried out in parallel as another part of the rendering pipeline. The bottom line is that greater realism in VR environments will not be obtained by pouring off-line CPU time and run-time space into rendering and recording characters in exacting detail; the visual effect of even the most perfectly animated figure is significantly reduced once the viewer recognizes that its movements are **exactly** the same each and every time it does something. We seek to capitalize on the intrinsically dynamic nature of interacting with and in a virtual world by recording less information and allowing motions to be modified on the fly to match the context in which they are replayed. Beginning efforts in this direction may be found in [14].

Real-time animation can be viewed as one of many "enabling" technologies for simulations. An animation, or visual simulation, of the activities and processes occurring in a simulation multiplies the effectiveness and communicative ability of the simulation, making its results more intuitively understood by non-expert viewers and participants. Unfortunately, when coupling real-time 3D animation to a simulation, one requires the presence of expensive rendering hardware in the simulation computer for generating the visuals. One also needs a high-performance general purpose processor, used for executing the simulation application itself, and for "feeding" the rendering hardware. Advances

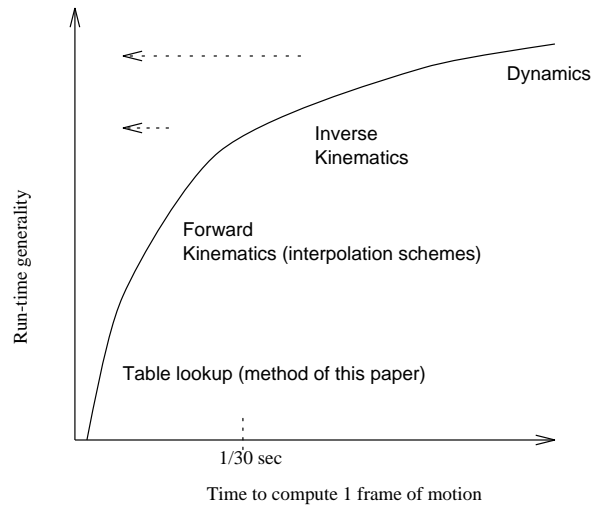


Figure 10: Intuitive relation between compute time and run-time generality for some motion generation techniques

in rendering hardware, as well as its general availability in low-cost PC platforms in the near future, may ameliorate this problem. Fast, 2 1/2 D animation techniques, such as those mentioned previously in discussing DOOM, have also found their way into most types of consumer interactive visual simulation applications (a.k.a. video games). This has heightened expectations, for quality animation on inexpensive platforms, from the users of these systems.

The general availability of mature simulation toolkits and libraries (e.g. SIMSCRIPT, Stella, and SimPack), on a variety of platforms, has greatly increased the use of simulation. However, there is a lack of such widely available, and used, toolkits for real-time animation. However, with the appearance of such toolkits as IRIS Performer, and OpenGL (an open 3D graphics rendering API, available on many platforms), it should become less burdensome in the future to add 3D animation to one's simulations.

8 Conclusions

We have described a system for off-line production and on-line playback of human figure motion for 3D visual simulation. The techniques employed are straightforward, and build upon several well known software systems and capabilities. As the number of possible states for a simulated human increases, the posture graphs will need to be replaced with a more procedural approach for changing posture. For applications built today on current workstations, the current technique is a balance between performance and realism.

Figure 10 shows a very coarse, and albeit intuitive, plot of the relation between run-time generality and computation time for several motion generation techniques. By “run-time generality”, we mean the notion of how general the types of motion are which can be generated by the algorithms. For example, table-lookup is very fast at run-time, but not very general. It can only generate the motion which has been recorded in its tables (of course, many things can be recorded and placed in the tables, but they are difficult to modify or “generalize” at run-time). On the other hand, a dynamics (or physical) simulation, with proper collision detection and response, can generate very realistic motion, under many different run-time conditions, so we deem it to be more “general”. But a full physical simulation is very expensive and time-consuming to compute, and some of the best for animation purposes are still at least an order-of-magnitude slower than real-time, for relatively simple environments [12]. For realistic agent animation in virtual environments, the research community will be trying to push this curve to the left and up, making the more powerful techniques run faster. The curve has been drifting to the left in recent years mainly on the progress made in rendering hardware and overall workstation compute performance. Interestingly, when an author’s motion generator can compute 1 frame of motion in less than 1/30th of a second, the author will usually claim it to be a “real-time” motion generator. Of course, our colleagues in the real-time community will start coughing very loudly at this. It is better to state that it is a “fast” motion generator, and leave it at that. The notions of “real-time” usually connote the presence of a scheduler of some sort, and some course of action in the event of a failure to execute by a specified deadline.

We chose humans for animating, as they are what we are interested in, but the techniques described in this paper could be applied to other complex articulated figures, whose states can be characterized by postures, and whose motions between postures can be organized into posture graphs.

References

- [1] Norman I. Badler, Rama Bindiganavale, John Granieri, Susanna Wei, and Xinmin Zhao. Posture Interpolation with Collision Avoidance. In *Proceedings of Computer Animation '94*, Geneva, Switzerland, May 1994. IEEE Computer Society Press.
- [2] Norman I. Badler, Cary B. Phillips, and Bonnie L. Webber. *Simulating Humans: Computer Graphics, Animation, and Control*. Oxford University Press, June 1993.
- [3] Jay Banchemo. Results to be published on system for dismounted infantry motion in a SIMNET image generator. Topographical Engineering Center, US Army.

- [4] Edwin H. Blake. A Metric for Computing Adaptive Detail in Animated Scenes using Object-Oriented Programming. In G. Marechal, editor, *Eurographics '87*, pages 295–307. North-Holland, August 1987.
- [5] David R. Pratt et al. Insertion of an Articulated Human into a Networked Virtual Environment. In *Proceedings of the 1994 AI, Simulation and Planning in High Autonomy Systems Conference*, University of Florida, Gainesville, 7–9 December 1994.
- [6] Paul A. Fishwick. *Simulation Model Design and Execution*. Prentice Hall, 1995. p. 114–135.
- [7] Thomas A. Funkhouser and Carlo H. Séquin. Adaptive Display Algorithm for Interactive Frame Rates During Visualization of Complex Virtual Environments. In James T. Kajiya, editor, *Computer Graphics (SIGGRAPH '93 Proceedings)*, volume 27, pages 247–254, August 1993.
- [8] John P. Granieri, Welton Becket, Barry Reich, Jonathan Crabtree, and Norman I. Badler. Behavioral Control for Real-Time Simulated Human Agents. In *ACM SIGGRAPH 1995 Symposium on Interactive 3D Graphics*, pages 173–180, 1995.
- [9] Institute for Simulation and Training, Orlando, FL. *Standard for Distributed Interactive Simulation - Application Protocols (v 2.0, 4th draft, revised)*, 1993.
- [10] Institute for Simulation and Training, Orlando, FL. *Enumeration and Bit-encoded Values for use with IEEE 1278.1 DIS - 1994*, ist-cr-93-46 edition, 1994.
- [11] Hyeongseok Ko. *Kinematic and Dynamic Techniques for Analyzing, Predicting, and Animating Human Locomotion*. PhD thesis, University of Pennsylvania, 1994.
- [12] Brian Mirtich and John Canny. Impulse-based Simulation of Rigid Bodies. In *ACM SIGGRAPH 1995 Symposium on Interactive 3D Graphics*, pages 181–188, 1995.
- [13] John Morrison. The VR-Link(tm) Networked Virtual Environment Software Infrastructure. *Presence*, 4(2):194–208, Spring 1995.
- [14] Ken Perlin. Danse interactif. SIGGRAPH Video Review, Vol. 101 1994.
- [15] Marc Raibert. Human Animation and Biomechanics. In *Proceedings of The First Workshop on Simulation and Interaction in Virtual Environments*, page 215, University of Iowa, Iowa City, IA, July 13-15 1995. (and also private communication).

- [16] Douglas A. Reece. Extending DIS for Individual Combatants. In *Proceedings of the 1994 AI, Simulation and Planning in High Autonomy Systems Conference*, University of Florida, Gainesville, 7–9 December 1994.
- [17] Douglas A. Reece. Soldier Agents in a Virtual Urban Battlefield. In *Proceedings of The First Workshop on Simulation and Interaction in Virtual Environments*, pages 204–209, University of Iowa, Iowa City, July 13–15 1995. (and additional private communication).
- [18] John Rohlf and James Helman. IRIS Performer: A High Performance Multiprocessing Toolkit for Real-Time 3D Graphics. In Andrew Glassner, editor, *Proceedings of SIGGRAPH '94 (Orlando, Florida, July 24–29, 1994)*, pages 381–395, July 1994.
- [19] Neil J. Rubenking. The DOOM Phenomenon. *PC Magazine*, 13(19):314–318, 1994.
- [20] Greg Turk. Re-tiling Polygonal Surfaces. In Edwin E. Catmull, editor, *Computer Graphics (SIGGRAPH '92 Proceedings)*, volume 26, pages 55–64, July 1992.
- [21] Frank Wysocki and David Fowlkes. Team Target Engagement Simulator Advanced Technology Demonstration. In *Proceedings of the Individual Combatant Modeling and Simulation Symposium*, pages 144–190, 15–17 February 1994. Held in Fort Benning, GA.
- [22] Jeffrey Adam Young. Doom's Day Afternoon. *Computer Player*, pages 20–28, October 1994.
- [23] Jianmin Zhao and Norman I. Badler. Inverse kinematics positioning using nonlinear programming for highly articulated figures. *ACM Transactions on Graphics*, to appear, 1995.