# Staying Oriented with Software Terrain Maps

Robert DeLine

Microsoft Research, Microsoft Corporation, Redmond, USA

## Abstract

*Developers often find themselves lost as they navigate around large programs, particularly when those programs are unfamiliar. This paper presents a new visualization, called a software terrain map, intended to keep a programmer oriented as she navigates around source code in the editor. The design is based on the metaphor of cartographic maps, which are continuous (no wasted space), have enough visual landmarks to allow the user to find her location perceptually rather than cognitively, and lend themselves to overlaying data. Although an optimal layout for software terrain maps is computationally intractable, the paper presents an efficient, heuristic algorithm that produces good results.*

## 1   Introduction

In our recent study in which experienced developers attempted to enhance an unfamiliar program, the participants consistently got lost in the source code. [1] They typically explored the code by opening and scanning many documents, by jumping back and forth among related definitions, and by iterating over the result sets of text searches. Despite their years of experience and their familiarity with the programming language, development environment, and problem domain, these navigation steps would quickly leave them disoriented. Many times, a participant would issue a new query to find a previously visited definition, and in a few cases, a participant would even inspect a previously visited definition and not recognize it.

Such disorientation is easy to understand. A typical program is very uniform in visual appearance and relies heavily on names to distinguish its parts. A typical development environment, like Visual Studio or Eclipse, also relies on names in its overview displays, like the tree of project files and the type hierarchy. Hence, staying oriented while navigating a program requires familiarity with its names, which places a large burden on both short- and long-term memory. In this paper, I describe a new display, intended to allow the programmer to use her spatial memory to stay oriented while navigating the source code.

To keep the user oriented during code navigation, I propose supplementing the development environment with an overview diagram to show the programmer's current location in the program ("you are here") and recent navigation steps (a vapor trail). In addition to reflecting the user's navigation in the editor, the overview would also allow navigation. For instance, clicking on the overview would also cause the editor to show the corresponding part of the code. The intent is to allow the programmer to use spatial memory to navigate to sought parts of the program. To realize this intent, several desiderata seem reasonable:

1.  The display should show the entire program. That is, whatever definition the programmer navigates to in the editor should have a representation at a reasonable level of detail in the display. Hence, the use of elision or abstraction to scale to large programs would not be appropriate. Elision would cause some navigation steps to be "off the map," while abstraction would cause a navigation step within an abstracted part of the program to appear as non-movement in the overview display.

2.  The display should contain enough visual landmarks to allow the developer to find parts of the program perceptually, rather than relying on names or other cognitive cues. For instance, to find Rome on a map of Europe, I scan for Italy's famous boot shape rather than reading for the word *Italy*. An overview display of software should have similar visual landmarks.

3.  The display should remain visually stable as the user navigates. If the display's content were to change as much as the editor's while the user navigates around the program text, it would provide little help in keeping the user oriented. Further, editing the program text should cause proportional changes to the display.

4.  Finally, to justify the additional screen space needed for the display, the display should be capable of showing global program information other than navigation steps. For instance, we might use the display to show program execution paths, like the call stack when an exception is raised or the hot path that a profiler reports, or to show team awareness data, like which developers are currently working on which parts of the program.

These desiderata mean that several popular technologies are not suitable for such a display. UML class diagram
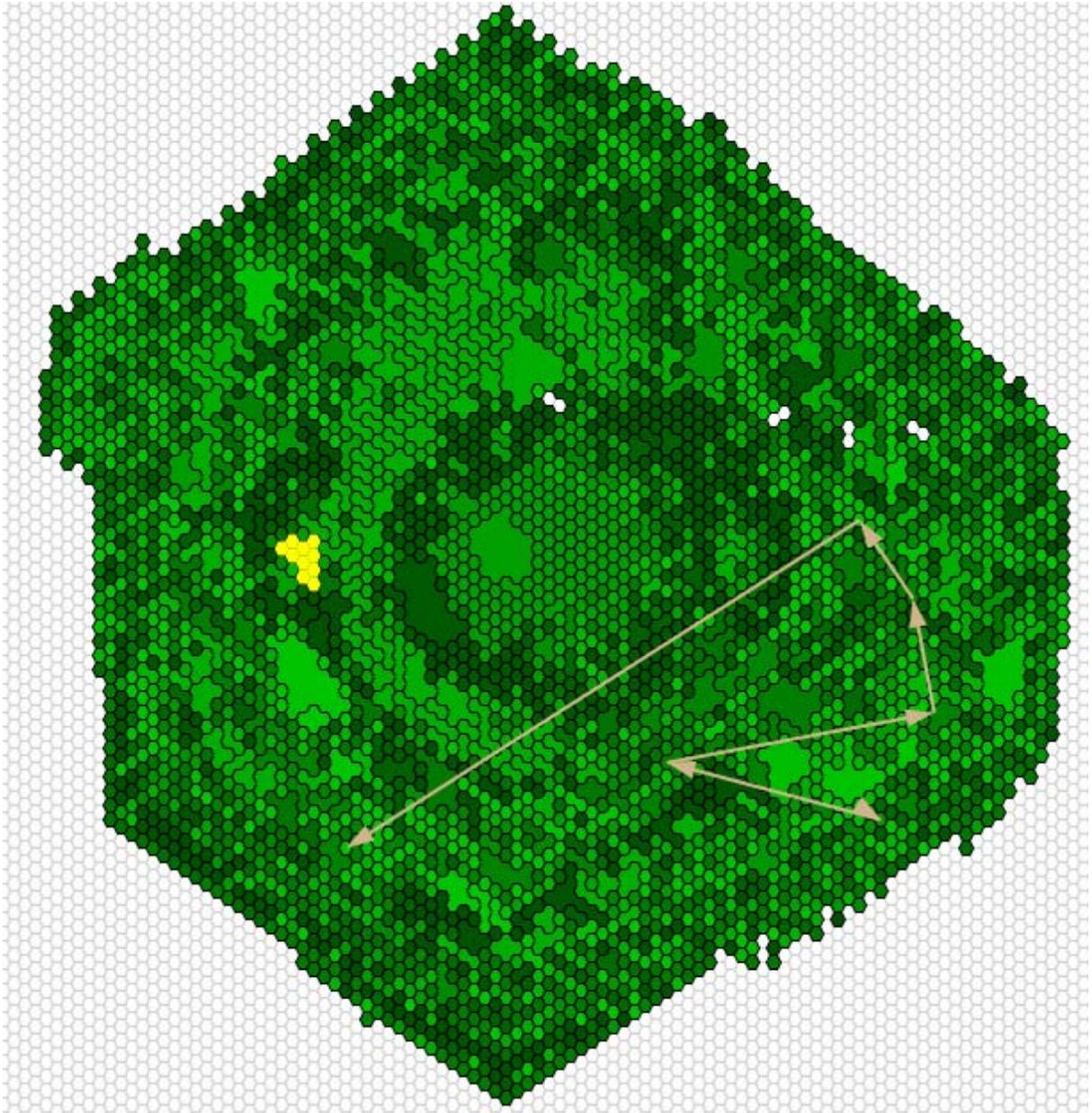
**Figure 1. The software terrain map of a library for analyzing .NET bytecode. Each of the 3800 regions is a method, whose size is proportional to the method's textual size (one tile for every two lines of code). Each class has its own shade of green. The highlighted (yellow) method is the one that the user is currently reading in the editor. The arrows show the program's execution path up to the current debugger breakpoint.**

and other box-and-line architectural diagrams, for instance, are not good candidates. First, they are often drawn at an inappropriately high level of abstraction or elide parts of the program to keep the diagrams small. A developer working on an object-oriented program navigates among and edits individual class members. Hence,

to provide a location marker while remaining visually stable, the display must show every member in the program all at once. This would be difficult to do with a UML class diagram. Second, such diagrams are visually uniform, particularly when they include many boxes. Distinguishing the boxes is more easily accomplished
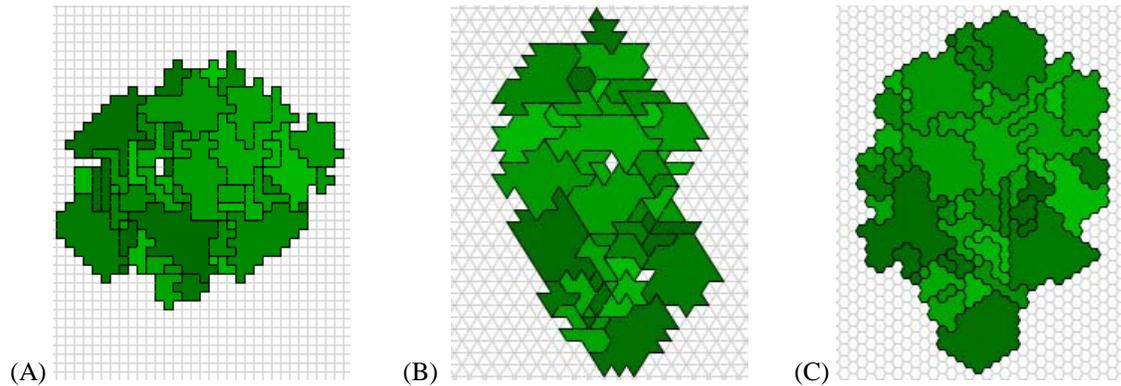
**Figure 2. Three software terrain maps for the same library, built on three grid systems: (A) squares; (B) equilateral triangles; and (C) hexagons.**

by reading their labels than by relying on visual cues like box position or the topology of edges between boxes.

A popular technique for scaling complex information to fit the screen is to use a detail-within-context display, like fisheye views.[2] Such a display gives more screen space and provides more detail about the user's current focus, diminishing and abstracting the other parts of the information. For instance, SHriMP uses this technique on box-and-line displays of software structure.[4] Using this technique for a navigation display would mean that the display would change every time the developer moved in the text editor. The use of animated transitions diminishes the disorientation caused by changing the display, but only when the user's attention is on the animation (and not on the editor, for instance). Even if the navigation display and editor occupy the same window, changing both simultaneously is confusing and unlikely to keep the user oriented.

## 2   Software Terrain Maps

To satisfy the desiderata above, I am designing a new software visualization called a *software terrain map,* based on the metaphor of cartographic maps. An example is shown in Figure 1. Cartographic maps have many nice properties: they are continuous (no wasted space) and stable (except perhaps at a geological time scale); they contain obvious and memorable visual landmarks (e.g. the shapes of boundaries, the positions of natural features); they lend themselves to overlaying data, both for easing navigation (e.g. names, roads, icons for features) and for conveying information in context (e.g. political, demographic, or economic patterns); and they are very familiar. A software terrain map is designed to show all of a software's parts, either behind the text in the editor window or on a second monitor. A highlight on the map continuously updates to show the part that the programmer is currently editing.

To mimic the continuous nature of cartographic maps, I partition the screen into tiles and assign tiles to the software parts. An algorithmically inexpensive approach is to choose locations for the parts and then to draw a Voronoi diagram around the locations to partition the screen. However, the shapes of the tiles constitute the display's major visual landmarks, and the use of Voronoi diagrams provides only indirect control over the tile shapes. Instead, I first partition the screen into regularly shaped tiles and then assign tiles to the software parts. As can be seen in Figure 1, the result gives the map an overall tidy, regular appearance, while containing enough irregularities to create visual landmarks.

Mathematicians, beginning with Golomb in the 1950s, have studied building shapes out of regularly tilings of the plane.[3] In particular, they have studied building shapes from squares, which they call *polyominoes* (see Figure 2A); from triangles, which they call *polyiamonds* (see Figure 2B); and from hexagons, which they call *polyhexes* (see Figure 2C). Software terrain maps can be drawn based on any of these three.

### 2.1   Layout Problem

To draw a software terrain map, we model the program as a set of *components* described by two metrics, which are parameters to the layout algorithm: Size($c$) which gives the number of contiguous tiles to assign to the component $c$; and Affinity($c_1$, $c_2$), which is the degree to which components $c_1$ and $c_2$ are related. The problem, then, is to locate components near each other in proportion to their affinity, while assigning each component the appropriate number of tiles. That is, computing the layout is an constrained optimization problem to find

min $\forall(c_1, c_2) \bullet$ Affinity($c_1$, $c_2$) $\times$ Distance($c_1$, $c_2$)

where Distance is a measure of screen distance between components, for example the Euclidean distance between their centroids.

This formulation of the layout problem is intentionally generic to allow me to explore various useful notions of size and affinity. For the terrain map in Figure
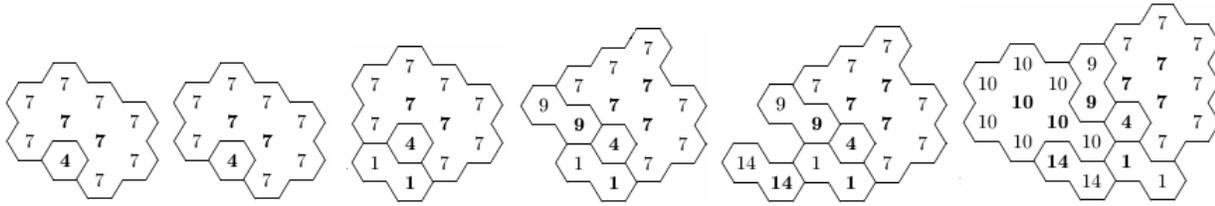
**Figure 3. The first six steps of the layout algorithm (left to right), showing the addition of components identified by numbers: component #7 has size 9; #4, size 1; #1, size 2; #9, size 2; #14, size 2; and #10, size 8. Each of a component's coordinates is labeled with the component's number; taken coordinates are in boldface.**

1, the components are methods and the number of tiles assigned to a method is proportional to its text size (in this case, one tile for every two lines of code). The intention is that methods that appear large in the editor also appear large in the terrain map, helping the programmer to learn the correspondence between the two.

For this figure, affinity is an arbitrary score chosen to reflect both control and data flow. In this case, to compute the affinity between $c_1$ and $c_2$, I give 20 points if $c_1$ calls $c_2$, 20 additional points if $c_2$ calls $c_1$, and one point for each field read or written in both $c_1$ and $c_2$. There is nothing "magical" about this formula. Part of the research agenda is to find formulas for affinity that have two nice properties: (1) the more often the developer navigates between two methods, the closer they appear on screen; (2) paths overlaid on the map (such as execution paths) are drawn as "snakes" rather than "scribbles." Ideally, a formula for affinity would involve information only about the program's static structure (since this is immediately available), but it may also be useful to include measures from the program's source history, from traces of its execution behavior, or even from the team's bug database or communications.

## 2.2 Layout Algorithm

The layout problem, while conveniently generic, is also computationally intractable, with clear relationships to both bin packing and the traveling salesman problem. Fortunately, there is a heuristic quadratic algorithm that produces appealing results. (Here, quadratic time is a lower bound for this problem since inspecting the N×N affinity matrix alone requires quadratic time. The bound could be lowered further by disregarding large portions of the affinity matrix.) The gist of the algorithm is to add each component to the map one at a time, spiraling outward from the center, and to adjust the shape of previously added components to improve their proximities to high-affinity components. This shape adjustment is limited to keep the overall algorithm quadratic.

The algorithm proceeds in two phases. First, we choose an order in which to add the components to the map. This order, perhaps even more than the size and affinity metrics, affects the map's final appearance, so I am experimenting with several choices. The order used to produce Figure 1 first sorts the methods by each method's total affinity for other methods. For each method $m$ in this sorted list, $m$ is added to the final ordering, then we add all those methods reachable from $m$ through a depth-first traversal of the program's call graph. At each stage of the depth-first traversal, the callees are explored in the order from highest to lowest total affinity. Two other approaches are to use a breadth-first search of the call graph and to ignore the call graph altogether and simply use the list sorted by total affinity. Of these three, the depth-first approach produces execution path overlays that are less "scribbly" than the other two approaches, at least in my initial experiments.

The second phase of the algorithm adds the components to the map in the order that the first phase determines. For each component to be added, we first compute the previously added component with the highest affinity for the new component, which I call the target. The second phase attempts to get the new component as close as possible to the target, without using more than linear time to find a good position. The second phase of the algorithm is parameterized by the grid to be used, namely, one of the three shown in Figure 2. A grid, as an abstract data type, supports a single operation: CoordinatesAtDistance($d$, ($x$,$y$)) returns the set of grid coordinates at Manhattan distance $d$ from coordinate ($x$,$y$). I'll use the expression Neighbors($x$,$y$) to mean CoordinatesAtDistance (1, ($x$,$y$)).

To allow a component's shape to be adjusted as new components are added, each component is assigned two types of coordinates: a component's *taken* coordinates are fixed (i.e. the taken coordinate belongs to the component now and forever); a component's *claimed* coordinates can be exchanged for other coordinates. The algorithm maintains the following invariants:

(1) the number of a component's taken and claimed coordinates equals the component's size, i.e. Taken($c$) + Claimed($c$) = Size($c$);

(2) the component's taken coordinates are all contiguous, i.e. $\forall (x,y) \in$ Taken($c$) • $\exists (x',y') \in$ Taken($c$) • $(x,y) \in$ Neighbors $(x',y')$; and

(3) every claimed coordinate is the neighbor of a taken coordinate, i.e. $\forall(x,y) \in$ Claimed$(c) \bullet \exists\ (x',y') \in$ Taken$(c) \bullet (x,y) \in$ Neighbors $(x',y')$.

Subject to these invariants, each component maintains as few taken coordinates as possible, since the more claimed coordinates a component has, the more flexible its shape, due to a process called claim renouncing, described below.

Pseudocode for the core of the algorithm is shown at right. For each component to be added, the algorithm begins looking at distance 1 from the target and proceeds to greater distances until enough room for the new component has been found. For a given distance, we first divide the coordinates at that distance into the empty ones (the ones that no component has claimed or taken) and the claimed ones. We first consider each empty space in turn as a possible root for adding the new component, turning to the claimed ones only if there are no suitable empty ones. At each root, we search for enough coordinates to assign to the component to make up its size. If we cannot find enough coordinates at that root, any state changes made to the grid are abandoned and we try the next candidate root.

To search for coordinates from a root coordinate, the component first takes the root coordinate. To find each additional coordinate needed, we search among the direct neighbors (coordinates at distance 1) of the component's taken and claimed coordinates. The component can grow to include a neighboring coordinate either if the coordinate is empty or if the coordinate is claimed by another component willing to renounce its claim (as described below). When the component finds a candidate neighboring coordinate, it claims it. If this newly claimed coordinate is a direct neighbor of one of the component's taken coordinates, the invariants are maintained and the algorithm can continue the search. However, if the newly claimed coordinate is a direct neighbor only of the component's claimed coordinates, then invariant (3) is violated. To re-establish the invariant, we convert one of the claimed coordinates to a taken coordinate and then continue the search.

For a component to renounce its claim on a coordinate, it must find a replacement coordinate to claim instead. The search for the replacement is exactly as described in the previous paragraph, with two exceptions. First, we must keep track of the coordinate being renounced so that the search for a replacement does not end up finding the one we want to renounce. In fact, since the search for a replacement can cause neighboring components to try to renounce their own claims, we must track all coordinates being renounced. (Otherwise, we can get cycles of neighboring components fruitlessly swapping renounced coordinates.) Finally, this recursive process of neighbors renouncing claimed is limited by a

```
type Component :
    var grid : Grid;
    var layout : GridLayout;

    PlaceNear (target : Component) :
        for distance := 1 to ∞ :
            candidates := layout.CoordsAtDistance(distance, target);
            foreach coord in EmptyCoordinates(candidates, grid) :
                if this.PlaceAt(coord) : return;
            foreach coord in ClaimedCoordinates(candidates, grid) :
                if this.PlaceAt(coord) : return;
    end PlaceNear

    PlaceAt (start: Coord) : bool
        grid.BeginTransaction();
        if grid.Claimed(start) ∧ ¬ grid.Claimant(start).Renounce(start,{},MAX)
            grid.AbortTransaction()
            return false;
        grid.Take(this, start);
        placesToExpand := new Queue<Coord>;
        placesToExpand.Enqueue(start);
        while ¬ this.CompletelyInGrid(grid) ∧ placesToExpand.Count > 0
            placeToExpand := placesToExpand.Dequeue();
            expanded := false;
            foreach c in grid.Neighbors(placeToExpand)
                if grid.IsEmpty(c) ∨
                    grid.IsClaimed(c) ∧ grid.Claimant(c).Renounce(c, {c}, MAX)
                    expanded := true;
                    grid.Claim(c, this);
                    placesToExpand.Enqueue(c);
            if expanded ∧ this = grid.Claimant(placeToExpand)
                grid.Take (placeToExpand, this);
        if this.CompletelyInGrid(grid)
            grid.CommitTransaction();
            return true;
        else
            grid.AbortTransaction();
            return false;
    end Place

    Renounce (toRenounce: Coord, forbidden: Set<Coord>, limit: int) : bool
        if limit = 0 : return false;
        othersClaims := {};
        foreach takenCoord in grid.TakenSet(this)
            foreach c in layout.Neighbors(takenCoord) \ forbidden
                if grid.IsEmpty(c)
                    grid.Unclaim(toRenounce, this);
                    grid.Claim(c, this);
                    return true;
                else if grid.IsClaimed(c) ∧ this ≠ grid.Claimant(c)
                    othersClaims.Add(c);
        foreach c in othersClaims
            if grid.Claimant(c).Renounce(c, forbidden ∪ {c}, limit-1)
                grid.Unclaim(toRenounce, this);
                grid.Claim(c, this);
                return true;
        foreach cl in grid.Claimed(this) \ forbidden
            foreach c in layout.Neighbors(cl) \ forbidden
                if grid.IsEmpty(c) ∨
                    grid.IsClaimed(c) ∧
                        grid.Claimant(c).Renounce(c, forbidden ∪ {c,cl}, limit-1)
                    grid.Take(cl, this);
                    grid.Unclaim(toRenounce, this);
                    grid.Claim(c, this);
                    return true;
        return false;
    end Renounce
end Component

type Terrain :
    SolveLayout (comps: List<Component>, layout: GridLayout) :
        target := new Map<Component,Component>;
        for i := 0 to size(comps)-1 :
            target[comps[i]] := CompWithMaxAffinity(Sublist(comps,0,i));
        grid := new GridState
        foreach comp in comps :
            comp.layout := layout;
            comp.grid := grid;
            comp.PlaceNear(target[comp]);
    end SolveLayout
end Terrain
```
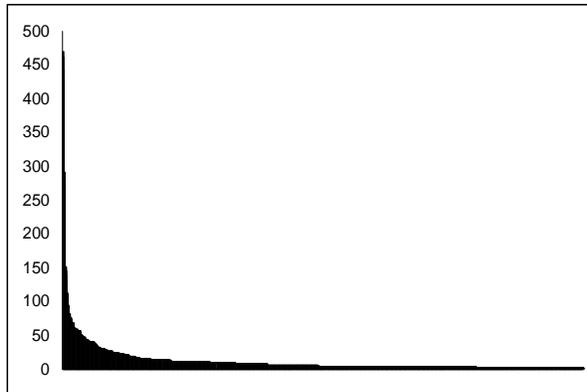
constant bound to ensure the question of whether a component may have a given coordinate can be answered in constant time.

Figure 3 illustrates this process with six components added to a grid. The components are identified by a number, and the a component's coordinates are labeled with its number. Claimed coordinates are shown in lightface; taken coordinates, in boldface. The first three components are added by taking empty coordinates. To add component #9 (of size 2), which has the highest affinity for component #4, component #7 renounces one of its claimed coordinates, so that #9 may have it. Similarly, in the sixth step, to add component #10, component #7 renounces a claimed coordinate to allow #9 to renounce one of its claimed coordinates so that #10 may have it. As the figure shows, the ability to renounce claimed coordinates allows components to get closer to their targets than they would if we were to use a pure greedy approach (all coordinates are taken).

## 3 Limitations and Next Steps

Although software terrain maps meet the desiderata mentioned in the introduction, there are drawbacks. First, basing the size of the methods on the size of the methods' text leads to many methods of size 1. The distribution of the sizes of this library's 3800 methods is an exponential decay curve, which is typical of several systems I measured:



The more size-1 methods there are in a software terrain map, the fewer visual landmarks.

One way to address the problem is to overlay additional landmarks on top of the methods, based, for instance, on the method's control structure. For example, one could add icons like the following, which are analogous to the symbols for schools, campgrounds, etc. found on cartographic maps:

| | |
|---|---|
| ◯ | loop |
| ◎ | nested loop |
| ☰ | switch statement |

These three are good candidates in that relatively few methods contain them. Such icons, however, are not useful at the scale of the map in Figure 1 since the individual tiles are too small to contain the icons.

A more serious limitation is that the appearance of the map is based on constraints which change as the software evolves. For instance, when the developer adds a new method, placing this method in the middle of the map would cause a cascade of claim renouncing, which would cause many methods to change shape. The result could be very disorienting. One approach to the problem is to keep all new methods off to the side, ignoring the new methods' affinities, until the developer (or perhaps the whole team) is ready for a large, disorienting map change. This solution does not apply to methods already in the map that are gaining new code. These growing methods would cause similar cascading changes. Experimenting with how visually disorienting users find these cascades is future work. Note that method deletions can be handled by leaving holes. The algorithm described here already generates a few holes (visible both in Figures 1 and 2), which I have allowed as another form of visual landmark.

I have implemented a prototype version of software terrain maps and have integrated it into Microsoft's Visual Studio development environment. The next step is a formal user study to evaluate how well these maps keep users oriented.

## 4 References

[1] R. DeLine, A. Khella, M. Czerwinski and G. Roberson, "Towards understanding programs through wear-based filtering," *Proc. Symp. on Software Visualization*, 2005.

[2] Furnas, G. "Generalized fisheye views," CHI '86.

[3] Golomb, S. *Polyominoes*, Princeton University Press, 1952.

[4] Storey, M.-A., "SHriMP views: an interactive environment for exploring multiple hierarchical views of a Java program," ICSE 2001 Workshop on Software Visualization, Toronto, Ontario, Canada, May 12, 2001.