

Analysis of a Metropolitan-Area Wireless Network

DIANE TANG

Stanford University, Gates 3A, Stanford, CA 94305-9030, USA

MARY BAKER

Stanford University, Gates 4A, Stanford, CA 94305-9040, USA

Abstract. We analyze a seven-week trace of the Metricom metropolitan-area packet radio wireless network to find how users take advantage of a mobile environment. Such understanding is critical for planning future large-scale mobile network infrastructures. Amongst other results, we find that users typically use the radios during the day and evening. Of the users who move around during the trace (over half), we find that the more locations a user visits on a daily basis, the closer together, on average, those locations are. While these results are only known to be valid for this particular network, we hope future analysis of other networks will add to a growing understanding of mobile network behavior.

Keywords: metropolitan-area wireless network, network analysis

1. Introduction

Currently, mobile and ad hoc networking is the topic of many research and development efforts. Much of this work focuses on providing future users with network resources and connectivity no matter where they are. Whether the work concentrates on adapting applications to changing user location or on devising new protocols to handle mobility, it is largely based on assumptions of how users will take advantage of a mobile environment. It is difficult to verify these assumptions since we are unaware of any publicly available studies of a sizeable metropolitan-area wireless network. Therefore, many research and development efforts must drive their simulations using assumed models of user movement not derived from observation.

In this paper, we analyze a network trace of the Metricom packet radio network, a metropolitan-area wireless network, to find answers to overall network questions such as when the mobile network is the most active, how active the network gets, where the network is active, as well as radio mobility questions such as how far, how often, and when users move. The answers to such questions are crucial in planning a future mobile network infrastructure, and in understanding how people actually take advantage of a mobile network. While these results are only known to be valid for this particular network, we hope future analysis of other networks will add to a growing understanding of mobile network behavior.

We present several results in this paper, including our finding that the more locations users visit on a daily basis, the closer together, on average, those locations are. In addition, the distance users move is a Gaussian distribution around the radius of the network. We also find that radios are used mostly during the day and evening hours.

In this paper, we first present some background information about the data before we present the actual results from the analysis. The analysis is divided into two parts: overall network behavior and radio mobility.

2. Background

In this section, we describe the network traced, how the data was collected, and some issues that arose in the data analysis.

2.1. Data collection

The traces we study here were obtained from Metricom [2,6]. Metricom has installed a RicochetTM packet radio network infrastructure in three major metropolitan areas (San Francisco Bay Area, Washington D.C., and Seattle), as well as in some airports, hotels, and college campuses scattered across the United States. This infrastructure consists of "poletop" repeaters distributed throughout the covered areas. Each poletop is one of two types:

- a wireless repeater, which just forwards packets on to another poletop via the radio interface, or
- a wired access point, which has both a radio interface and a wired interface. Typically, wired access points have 8, 16, or 24 radios on them, as they are the focal point of many other wireless repeaters.

The range of a repeater is normally about half a mile. This range may vary depending on external conditions, such as the weather or the location of buildings and hills.

When a subscriber radio is first turned on, it scans the network for poletops, and chooses one with which to register. This choice is usually based on signal strength, but may also be based on load-balancing considerations. This poletop is responsible for forwarding all packets to and from the Metricom network on behalf of that radio. Radio registrations also occur whenever the radio changes its primary poletop and according to a predetermined pattern when the radio is stationary (see section 5.2). Note that while a radio registers with only one poletop at a time, the radio does keep an internal list of all poletops within its range. Each registration, whether at the same or at a different primary poletop, is logged by a central nameserver.

The trace consists of a nameserver log covering a sevenweek period from February 1, 1998, through March 23, 1998, with the exception of three holes in the trace data, during which no registrations were logged: February 16, 6 a.m. through 1 p.m.; February 16, 5 p.m., through February 17, 4 a.m.; and February 17, 4 p.m., through February 18, noon. The network at the time of the trace consisted of 14,053 poletops and 24,773 radios.

There were a total of 7,726,678 events logged over these seven weeks. Of those, 5,982,846 are registrations and the other 1,743,832 are queries. A registration occurs when the radio informs the nameserver of its current primary poletop, i.e., the poletop to which packets destined for this radio should be sent. A query occurs when the radio queries the nameserver about some other entity in the network, such as another radio or poletop. Queries are usually made at the start of a connection, and the radio can register at different poletops while a connection is on-going.

Each log entry consists of the following information:

- a timestamp taken at the nameserver with accuracy to the second,
- a sequence number,
- the radio id,
- the wired access point used,
- the name of the radio's choice for primary poletop,
- whether the entry is a registration or a query.

A poletop's name is an encoding of its latitude and longitude, so given a poletop's name, we can determine its geographic location and therefore the approximate geographic location of the radio registering with it.

2.2. Data analysis techniques

In analyzing this trace, the main difficulty is in dealing with the sheer volume of data -7,726,678 events, 24,773 radios, and 14,053 poletops – making it impractical, if not impossible, to look at each radio by hand. As a result, we needed methods to help gather our results automatically. We turned to a technique commonly used in data mining and machine learning called clustering, also known as an unsupervised learning technique.

Clustering algorithms take a set of points in n-dimensional space and find coherent subsets. Each subset consists of points that are clustered together. The advantages of using clustering algorithms are the ability to categorize radios automatically, and the ability to find groupings that we might not otherwise find. The disadvantages are that the results are dependent on the parameters and distance functions used.

We use three different clustering algorithms:

- *k*-means [9], an iterative clustering algorithm,
- hierarchical agglomerative clustering [8], a tree-forming clustering algorithm, and
- expectation-maximization (EM) [3], a method to find the means and variances of a mixture of Gaussian distributions. EM is especially useful when the range of values differs widely between dimensions.

We use both hierarchical agglomerative and k-means clustering algorithms to help determine when a radio has moved by grouping poletops into clusters, which we also refer to as locations. Without physically moving, the radio may still be within range of up to 20 different poletops at the same time and may, therefore, register at any of these poletops. We determine that a radio has moved when it registers with a poletop in a different cluster (a different location). In section 4, we use both EM and k-means to categorize the radios into different patterns of mobility.

3. Overall network behavior

The first three questions we asked have to do with overall network usage:

- 1. When is the network the most active?
- 2. How active does the network get?
- 3. Where is the network the most active?

The answers to these questions can help network planners plan future extensions to the network infrastructure. Also, they provide a way both to compare a simulated network to an actual network and to provide a basis for simulated networks.

To summarize the results presented in this section, we find that the network is more active on weekdays than on weekends, that the most active poletops handle up to 182 distinct radios an hour and handle over 1,600 registrations an hour, and that the network is most active where there is a high concentration of technical people. We now look at each question in more detail.

To answer the first question, figure 1 shows that the network is more active on weekdays, especially the days in the middle of the week. This pattern holds regardless of whether we are looking at all events, registrations only, or queries only. Figure 2 shows that the network is least active between 3 a.m. and 5 a.m., when most people are asleep. The two lines in figure 2, one lower in the evening hours and one higher, correspond to weekend and weekdays respectively. Even in Silicon Valley, where the majority of radio users and poletops are, people are less likely to work on weekends. There is a slight decrease in the number of active radios during weekday evenings, perhaps corresponding to when people stop working for the day. There is also a slight rise around 8 p.m. on weekends. We define an active radio as one that has at least one registration or query logged at the nameserver within the given time period.



Figure 1. Histogram of the number of active radios on each day of the trace. One corresponds to Sunday, two to Monday, etc. There is a dip around the 16th through the 18th days corresponding to the holes in the trace data.



Figure 2. Graph showing the average, minimum, and maximum number of active radios for each hour of the trace. The two sets of lines, one lower in the evening hours and one higher, correspond to the weekend and weekdays respectively. We ignore those days containing the holes in the trace during which no events were logged.

To find out how active the network gets, for each poletop we first count how many distinct radios register with or query from that poletop over the course of the trace. This ranges from handling only one distinct radio event over the course of the entire trace to handling 6,677 events (6,064 registrations, 434 queries), which occurs at the poletop at Metricom headquarters. Figure 3 shows the distribution within this range.



Figure 3. Cumulative histogram of the number of poletops (*y*-axis) that handle a certain number of distinct active radios over the course of the entire trace (*x*-axis). Note that to show the detail, we cut off the tail of the graph, which extends out to 6,677 radios for all events, 6,064 radios for registrations only, and 434 radios for queries only.



Figure 4. Two-variable histogram, where the darkness of the bar reflects the number of poletops that receive events from a given number of radios, by hour of the day. No matter what time of day it is, most poletops see events from fewer than 10 radios in an hour. The largest number of radios a poletop receives events from within an hour is 182. Note that the maximum number of radios at any point in the graph is the total number of poletops (14,053) multiplied by the number of times that hour occurs during the trace (once for each of the 52 days).

We also examine how many radios and radio events per hour a poletop needs to handle. Figure 4 shows that while the majority of poletops only handle one radio an hour, some



Figure 5. Two-variable histogram, where the darkness of the bar reflects the number of poletops receiving a given number of events per hour, by hour of the day. Note that the top of the graph, which extends to 1,601 events, has been cut off to show the detail.

poletops may handle up to 182 distinct radios an hour. Figure 5 shows that most of the time, most poletops handle 500 or fewer radio events per hour, although some particularly busy poletops handle up to 1,600 events in an hour. Note that the shapes of both the dark areas and the peak values in figure 4 resemble figure 2 with the low points falling between 3 a.m. and 5 a.m., whereas in figure 5, only the dark areas at the bottom portion of the graph resemble figure 2. The peaks in figure 5 do not follow the rise and fall of radio activity. This lack of a clear rise and fall in peak activity could be due to frequent radio registrations regardless of radio activity.

We also look at how often the network needs to handle events, both in the overall network and per poletop. In the overall network, the three largest intervals between two successive events by any radio, discounting the holes in the trace, are 1293 s (21 min), 744 s (12 min), and 59 s. Each of these intervals occurs only once. Figure 6 shows that by far, however, the network had either no time, or one to two seconds between events. If we look only at registrations in the network, the distribution of intervals is almost identical. However, if we look only at queries, we see that the curve is shallower, with some larger intervals up to 10 to 12 s occurring more frequently. This result is due to the many fewer queries than registrations that occur.

Figure 7 shows the cumulative distribution of intervals between events at distinct poletops, rather than in the overall network. This figure differs significantly from figure 6, with a much longer tail, due to different poletops having different loads. Some poletops are basically inactive, processing only two or three events during the entire trace, while other pole-



Figure 6. Cumulative histogram of the interval between successive events over the entire trace as a function of how often that interval occurs. For clarity, we do not plot the tail, which extends to 1,293 s.



Figure 7. Cumulative histogram of the interval between successive events at a poletop as a function of how often that interval occurs. For clarity, we do not plot the tail of the graph, which extends out to 4,319,696 s (71,994.9 min, 1199.9 hours, or 49.9 days) for all events and registrations only, and to 4,370,269 s (72,837.8 min, 1213.9 hours, or 50.6 days) for queries only.

tops handle thousands in an hour. One similarity between the two graphs is the lower curve for queries only, in both cases due to there being fewer queries overall.

Finally, figure 8 shows the most active poletops in the San Francisco Bay Area. A picture of the entire United States, covering all areas in which Metricom has established an infrastructure, is too crowded, so we choose to show just the



Figure 8. Picture of the San Francisco Bay Area. Each dot corresponds to a poletop. The darker the dot, the more radios visit that poletop over the course of the trace.



Figure 9. Histogram of the distance between successive events in the network by any radio as a function of how often that distance occurs.

Bay Area since it has the highest concentration of poletops. The darker dots correspond to the more active poletops. Several hot spots are worth noting. First, the dark area in San Francisco corresponds to the Financial District. The second strip of scattered dark dots running northwest-southeast corresponds to Highway 101, on which many high tech companies have their headquarters. Another hot spot is the dark dot in the upper part, which corresponds to Berkeley. Also, towards the bottom is a hot spot corresponding to Metricom headquarters. The lowest island of activity corresponds to Santa Cruz.

We also want to investigate how the activity in the network is distributed geographically, rather than temporally as shown in figure 6. Now we look at the interval in distance between successive events in the entire network. As we can see in figure 9, there are several common distances: 100 miles or less, around 650 miles, around 1000 miles, 1400 miles, 2000 miles, and 2400 miles. These distances all correspond either to traffic within a Metricom installation or to the approximate distances between various Metricom installations. Because the Bay Area contains about 65% of the poletops in the Metricom infrastructure, distances of 100 miles or less are prevalent. In other words, two successive registrations in the network are more likely to both be in the Bay Area than to be distributed across the country. However, we can see a Gaussian distribution around each of the common distances. This implies that given the radius of the network, the distance between registrations of any radio in the network is likely to be a Gaussian distribution around the radius.

4. Radio mobility

The other major set of questions we want to answer concerns radio mobility:

- 1. How often do radios move?
- 2. How far do they move?
- 3. Can we identify patterns of mobility?

Answering these questions is crucial for understanding whether and how people actually take advantage of a mobile environment. Also, understanding current radio mobility helps in choosing parameters for simulations of mobile networks. Unrealistic movement models lead to unrealistic simulation results.

To summarize the results in this section, we first find that users who are mobile do not move frequently, and that 64% of all users only appear at one location per day. As for how far users move, most users move within their local area, with fewer users traveling the long distances between different Metricom installations. In addition, as the number of locations visited by a user increases, the average distance traveled between each location decreases. Finally, we are able to find patterns of mobility for users, such as the number of users who are active both day and night versus users active only during the day. We now examine these findings in more detail.

We first calculate the number of different poletops and locations at which radios register over the course of the trace. Figure 10 shows that 42% of all radios are stationary, and that 64% of all radios visit only one location a day. This implies that although users do move around with their laptops, the movement is not very frequent. We can also see the difference between looking at poletops versus locations. While 42% of radios are stationary with respect to location, only 16% of radios are stationary with respect to poletops, meaning that radios often register with poletops that fall in the same cluster or general area. Our location finder underestimates user



Figure 10. Cumulative histogram of the percentage of radios that visit a given number of poletops, locations, or locations per day, over the course of the trace. For example, 64% of radios visit only one location per day, 42% visit only one location over the course of the entire trace, and 16% of radios visit only one poletop over the course of the entire trace. Note that to show the detail, we cut off the tail of the graph, which extends to 25 locations per day, 176 locations, and 1,025 poletops.



Figure 11. Histogram of the average, minimum, and maximum number of radios changing locations during a particular hour of the day. The time plotted is the second of two successive events.

mobility, so there may be more actual locations than we have found (see section 5.3).

We also measure how many radios appear, i.e., register or query, at a different location from their last appearance. Figure 11 shows these results. Unsurprisingly, this rate mimics the usage graph in figure 2, with anywhere from 122 (0.5%)



Figure 12. Cumulative histogram of the number of events as a function of the distance traveled. Note that the tail of the graph (which extends to 2,576 miles) has been removed to show detail. Also note that each pair of corresponding poletop and location lines converge at one mile. The straight line between 0 and 1 miles for the location graphs reflects the fact that locations cluster poletops together: poletops that are close together are considered to be a single location, and therefore, represent a movement of distance zero.

to 1,484 (6.0%) radios changing location at a time. There are, however, some interesting features to the weekday curve in figure 11. First, there is a dip around 11 a.m. to 2 p.m., approximately corresponding to lunchtime, showing that many people do not take their laptops to lunch wih them. Second, there is an increase in the number of radios changing location around 5 p.m., about when many people stop working for the day. We still observe the dip in the weekday evenings similar to figure 2.

We next look at how far a radio moves when it changes its primary poletop or location. Figure 12 shows the distribution of how often radios move a given distance, in terms of movement both between poletops and between locations. Over 79% of the events (or 80% of registrations, 90% of queries) involve no change in location. If we then look at figure 13, we can see that while the majority of successive events occur within 50 miles of one another, there are small peaks corresponding to the long-distance routes between airports and major cities in which Metricom has established an infrastructure. While this graph looks similar to figure 9, cross-country registrations are much less frequent when examining individual user behavior rather than overall network behavior.

We further examine radio movement in table 1, which shows the breakdown of radios by how often they move, both overall and per day, and how far they move. First, a radio that visits more locations overall does not necessarily also appear at proportionally more locations per day. For instance, the majority of radios that visit three or more locations overall still only appear at one to two locations on a daily basis. Furthermore, only 3.5% of radios that visit more than four locations overall appear at more than four locations per day on average. (Note that one distinct location may be counted as multiple visited locations per day. For example, a user that visits location A followed by location B before returning to location A in one day is counted as appearing at three locations in that day.) This result implies that users who visit many locations do not necessarily move around more on a daily basis,



Figure 13. Histogram of the percentage of successive events from a radio as a function of the distance traveled between those events. Note the logarithmic scale on the y-axis.

but rather visit the different locations over a longer period of several days or weeks.

More surprisingly, as the number of locations visited on a daily basis increases, the average distance traveled decreases. In other words, the more someone moves around each day, the closer the locations are. This overall decreasing trend may correspond to the flatter distribution of the range of distances among users who visit multiple locations per day, although the majority of locations are still relatively close together. This distribution could stem, for example, from users using their laptops on a train, such as the train that runs along Highway 101, causing them to "visit" numerous locations on their way to work.

Users who visit one location per day, independent of the total number of locations visited throughout the trace, have a large discrepancy between the average distance traveled and the median distance traveled. This discrepancy is likely the result of these users actually being split into two sub-groups: those who travel locally (for example, people who take their laptops home on the weekends), and those who travel long distances (perhaps people who only use their laptops when travelling). The median distance traveled reflects the usage patterns of this first sub-group, while this second sub-group of users increases the average distance traveled for all the groups of users who, independent of the total number of locations visited throughout the trace, only visit one location per day.

Despite this overall decreasing trend in the average distance, the median distance stays relatively constant, reflecting the many users who only move short distances. Further, the

 Table 1

 How far users move depending on how many locations and locations per day a user visits. For example, 510 users visit only two locations throughout the entire trace, but on average, visit between one and two locations per day.

Number of locations	Number of radios (%)	Average distance	Median distance	Number of locations per active day (max)	Number of radios (%)	Average distance	Median distance
1	10,459 (42.2%)	0	0	1 (1.0)	10,459 (100%)	0	0
2	2,746 (15.1%)	6.9	2.9	1 (1.0)	2,152 (84.1%)	11.6	1.6
				1-2(1.98)	510 (13.6%)	2.1	1.7
				2 (2.0)	84 (2.3%)	1.7	1.3
3	2,371 (9.6%)	5.2	1.6	1 (1.0)	1,367 (57.7%)	11.3	1.7
				1-2(1.98)	914 (38.5%)	3.5	1.6
				2 (2.0)	47 (2.0%)	2.3	1.9
				2-3 (2.92)	31 (1.3%)	2.3	2.0
				3 (3.0)	12 (0.5%)	2.2	2.2
4	1,694 (6.8%)	4.2	1.5	1 (1.0)	569 (33.6%)	8.4	1.6
				1-2 (1.98)	1,013 (59.8%)	4.6	1.5
				2 (2.0)	44 (2.6%)	2.3	1.5
				2-3 (2.98)	55 (3.2%)	2.1	1.5
				3 (3.0)	3 (0.2%)	2.0	1.9
				3-4 (3.65)	8 (0.5%)	1.5	1.3
				4 (4.0)	2 (0.1%)	7.9	2.7
>4	6,503 (26.3%)	7.9	2.4	1 (1.0)	438 (6.7%)	12.2	1.7
				1-2 (1.98)	4,434 (68.2%)	11.2	1.8
				2 (2.0)	126 (1.9%)	4.2	1.3
				2-3 (2.98)	1,020 (15.8%)	4.9	1.7
				3 (3.0)	36 (0.5%)	6.0	1.7
				3-4 (3.96)	210 (3.2%)	4.3	2.0
				4 (4.0)	13 (0.2%)	4.3	2.0
				>4 (24.48)	226 (3.5%)	8.0	6.9





Figure 14. Cumulative histogram of the number of times a certain number of registrations occur before a query. Note that to show the detail, we cut off the tail of the graph which extends to 7,948 registrations between queries.

average and median distances, in terms of the total number of locations visited rather than the average number of locations visited per day, are also relatively constant, implying that movement is dependent more on how many locations a user visits on a daily basis than on how many locations a user visits overall.

4.1. Radio queries

Before we look into categorizing radios into different patterns of mobility, we examine radio behavior with respect to queries. Queries are usually made when a radio makes a new connection, although some radios are used as part of Metricom's diagnostic utilities and make many more queries than are normal. Note that a query is made regardless of whether the user is using the radio in modem mode or StarMode (or packet mode) [2]. If the radio is being used in StarMode, the radio could have multiple outstanding connections at a time. Since we do not have information about the end of a connection, we do not account for this situation. Using the radios in modem mode is prevalent, however. Because very few radios are used for diagnostic purposes and very few radios are used in StarMode, we believe that our results are a good approximation.

Our main interest with regard to queries is in how much activity there is between queries. We first look at how many times radios register between queries. Figure 14 shows the cumulative distribution of how many registrations occur before a query, which is indicative of the amount of activity between the beginnings of different connections. The majority of connections occur after 30 or fewer registrations. In figure 15, we examine how often a radio changes its location or poletop before starting a connection. This graph is indicative of how

Figure 15. Cumulative histogram of how often a certain number of poletops or locations are registered at before a query is made. Note that to show the detail, we cut off the tail of the graph which extends to 1,193 locations and 3,334 poletops.

much movement there is between connections. Note that a connection can be long-lived and last for days, explaining the tail of the graph, which extends out to 1,193 locations or 3,334 poletops. Also note that we count the number of changes in location; if a radio visits location A then B then C then B and then A again, we count this sequence as four location changes. We see in figure 15 that 70.4% of connections have no movement before them at all, while 95% of radios change location fewer than four times. In other words, most radios do not move around much between connections. Unfortunately, we do not have information about when connections terminate so we do not know how much of the movement occurs while the connection itself is maintained.

4.2. Patterns of mobility

Our final goal is to categorize the radios into different patterns of mobility. One expected pattern, for example, might correspond to a user who moves between using the radio at home and at the office. Another pattern would correspond to the user who just uses the radio at home. Yet another pattern might correspond to a traveling salesman who uses the radio to keep in touch with the warehouse. We are interested in these patterns to see whether and how people take advantage of mobility.

Without comprehensively going through the radios by hand, there is no way to find all the categories, much less categorize each radio, without using clustering algorithms. We use a two-stage clustering because it avoids the difficulty in finding a single stage clustering that balances both the location and time information together. By separating out the clustering into two stages, we can first cluster the radios based on how mobile they are and on some overall frequency parameters. Once we know how mobile the radios are, we can adjust the second stage to reflect the results of the first stage, and sub-cluster based on time of day to determine when radios are active.

More specifically, in the first stage, we cluster the radios based on mobility using nine parameters:

- the total number of unique locations a radio visits throughout the trace,
- the total number of unique locations a radio queries from throughout the trace,
- the average number of locations a radio visits per active day,
- the average number of locations a radio queries from per active day,
- the total number of active days for a radio,
- the total number of query days for a radio,
- the average number of active days in an active week for a radio,
- the average number of query days in an active week for a radio, and
- the average number of query days per week for a radio.

An active day is defined as a day during which the radio visits some poletop, a query day is defined as a day during which the radio makes at least one query, and an active week is defined as a week during which at least one day is active. We use EM (expectation maximization) to do this first-stage clustering, since EM takes the differing ranges of the parameters into account, whereas the other clustering algorithms, such as k-means, do not.

We then take each cluster resulting from the first stage, and sub-cluster it based on more specific time information, namely the times of day that the radio is used. Specifically, for each radio we create two histograms counting the number of times the radio has an event per hour of the day, one for weekdays and one for weekends. These histograms show when the radio is active. We then normalize these histograms and use them as the parameters for the second-stage clustering, thus grouping radios with similar active times. We do this sub-clustering, rather than just re-cluster the entire space of radios, so that we can see how each mobility pattern changes by time of day.

We now look at the results from the first-stage clustering and second-stage sub-clustering, shown in tables 2 and 3, respectively. The first-stage clustering results in eleven clusters. Cluster 1 consists of radios that are essentially not used during the trace. Cluster 2 is comprised of stationary radios, perhaps the ones used on desktop machines for home connectivity. Cluster 3 consists of radios that are minimally mobile, i.e., users that predominantly use their radios from a single location, but may occasionally travel to the office to download some software or data. Clusters 4, 5, 6, and 7 all consist of radios that, while they do visit several locations, do not neces-

 Table 2

 Results from the first-stage clustering, which results in 11 clusters. Each cluster represents a parameterization for a different pattern of mobility.

Cluster	Number of radios	Total number of loc's (all, σ , qry only, σ)	Avg number of loc changes per active day (all, σ , qry only, σ)	Total number of active days (all, σ , qry only, σ)	Avg number of active days per active week (all, σ , qry only, σ)	Avg number of queries per day (σ)	Description
1	5,877	1.02 (0.17) 0 (0)	1 (0) 0 (0)	2.98 (7.34) 0 (0)	1.41 (0.95) 0 (0)	0 (0)	Not used
2	4,608	1 (0) 1 (0)	1 (0) 1 (0)	27.65 (18.43) 20.31 (15.69)	3.98 (1.94) 3.12 (1.68)	2.41 (2.21)	Stationary (home connectivity?)
3	4,765	2.98 (1.50) 1.60 (0.81)	1.40 (0.69) 1 (0)	22.66 (17.02) 16.04 (13.84)	3.51 (1.79) 2.77 (1.49)	2.28 (1.51)	Minimal mobility, moderate usage
4	1,117	5.59 (3.67) 2.87 (2.12)	2.81 (1.66) 1.26 (0.77)	5.63 (3.21) 3.78 (2.85)	1.76 (0.60) 1.38 (0.86)	3.91 (4.32)	Moderate mobility, minimal usage
5	1,422	5.13 (2.53) 3.65 (1.57)	1.50 (0.41) 1.16 (0.10)	15.50 (5.75) 13.63 (5.17)	2.66 (0.65) 2.51 (0.63)	2.42 (1.05)	Moderate mobility, low usage
6	1,603	4.55 (2.32) 3.17 (1.26)	1.39 (0.36) 1.12 (0.10)	35.22 (6.88) 31.23 (6.26)	4.69 (0.73) 4.24 (0.70)	2.59 (1.23)	Moderate mobility, moderate usage
7	984	5.19 (3.0) 3.43 (1.66)	1.84 (0.79) 1.25 (0.26)	47.99 (2.77) 46.16 (3.22)	6.05 (0.33) 5.83 (0.39)	3.64 (2.04)	Moderate mobility, high usage
8	882	5.23 (3.98) 2.41 (1.12)	9.43 (8.14) 1.15 (0.17)	37.62 (15.36) 17.41 (13.19)	5.15 (1.48) 2.76 (1.39)	1.68 (0.66)	Moderate mobility, train users?
9	1,459	12.05 (6.44) 6.85 (3.46)	3.06 (1.49) 1.43 (0.28)	22.21 (7.26) 19.67 (6.88)	3.51 (0.78) 3.34 (0.81)	3.23 (1.80)	High mobility, moderate-to-low usage
10	1,299	12.14 (7.83) 6.83 (4.09)	4.17 (2.51) 1.55 (0.39)	42.91 (5.63) 39.50 (5.88)	5.51 (0.62) 5.13 (0.66)	3.18 (1.33)	High mobility, high usage
11	757	18.85 (19.70) 9.41 (10.57)	6.87 (6.02) 1.80 (1.31)	30.92 (15.08) 21.34 (16.80)	4.54 (1.47) 3.33 (2.07)	12.45 (21.83)	High mobility, moderate usage, high connection rate

Table 3

Results from the second-stage clustering. Each row corresponds to a time pattern, which is shown with a representative histogram. The *x*-axis of the histogram is the hour of the day, and the *y*-axis is the percentage of events occurring during that hour. The black bars correspond to weekday usage, and the white bars correspond to weekend usage. Eight time patterns emerge, and the distribution of these patterns differs amongst the mobility patterns resulting from the first-stage clustering. Note that each first-stage cluster corresponds to one column of the table.



sarily visit all of those locations during a single day. In other words, the movement is spread out over a larger amount of time. The main difference between these three clusters is in how often the radios are used, with the radios in cluster 4 used the least, and the radios in cluster 6 used the most. However, the radios in cluster 4 also move around the most per day of usage, with more location changes per day. Cluster 8 consists of radios that on any given day, visit approximately twice as many (not necessarily distinct) locations as they see different locations in the entire trace. This behavior could correspond to users who commute by train and use their laptops while on the train, thus visiting most of their locations twice in one day, once on the way to work and once on the way back. Clusters 9–11 correspond to radios that move around much more, differing mainly in how often the radios are used and how many connections are formed per day.

Table 3 shows the eight prevalent timing patterns resulting from the second-stage clustering:

- 1. A smooth curve similar to figure 2, with usage balanced across weekends and weekdays.
- 2. A curve similar to the previous one, but the times are shifted so that the time of least usage is later (around 5 to 6 a.m. rather than 4 a.m.), perhaps consisting of techies who are used to working later hours.
- A distribution in which the usage later in the day primarily occurs during the weekend and the usage during the day primarily occurs on weekdays.
- 4. A distribution in which there is almost no activity all night, and the peaks are a little more obvious. In fact, slight peaks around 9 to 10 a.m. and 6 to 8 p.m. can be seen in the weekday and weekend usage.
- 5. A distribution with very obvious peaks, usually one to two big weekend peaks and two to three smoother weekday peaks. The histogram pictured is merely representative, since other time usage distributions in this cluster have the peaks at different times.
- 6. A distribution we call system administration hours, with usage throughout the day, but higher usage late at night.
- 7. A relatively smooth weekday usage curve, predominantly during the day, with almost no weekend usage.
- 8. And finally, a very sporadic distribution, with one to two peaks during both the weekday and weekends. Similar to 6, these peaks can be shifted in time.

From table 3, we see that while we can observe patterns of mobility derived from the first-stage clustering, not all users within a cluster execute their particular mobility pattern at the same time. Instead, users can be further divided into time patterns. For example, cluster 2 consists of 4,608 stationary radios. This set of radios can further be broken down into users who balance their radio usage across weekends and weekdays (time pattern 1), and users who almost never use their radios at night (time pattern 4, and to a lesser extent, time pattern 3), for example. Also note that not all time patterns occur in every mobility pattern and that the percentages of different time patterns occurring in different mobility patterns vary.

5. Lessons learned

During the course of this analysis, we learned several lessons about the importance of visualization tools, the difficulty in analyzing a trace taken for a purpose other than our own, and the importance of parameter choice in clustering algorithms.

5.1. Visualization

The first lesson we learned is that visualization tools are crucial given the volume of data. Without the ability to see a high-level view of the data, it is easy to get bogged down in the details of a very small part of the data set.

In addition, without using different levels of detail, visualization is very difficult, especially given the volume of data involved. In other words, rather than always dealing with the raw registration data, determining intermediate parameters, such as locations and the number of active days, is critical in being able to understand the overall picture.

We used a custom visualization built on top of the Rivet visualization environment [1], combined with our own implementation of the clustering algorithms to help us understand (and debug) the clustering algorithms themselves and to understand the results from the clustering algorithms. Having a custom visualization in which we could implement the exact visualizations and interactions we wanted, adjust the algorithm parameters, and implement and adjust the algorithms themselves was key in understanding the data and the results.

5.2. Tracing

The next lesson we learned was about the difficulty in analyzing a trace that was not gathered for the purpose of our analysis, since the trace gathered by Metricom for their own purposes is missing data we would have found useful. Several unanswered questions and ambiguities are a direct result of not having control over what data was recorded and where the data was recorded.

First, one question we wanted to answer was how long a radio stays active, or in other words, how long sessions last. However, session delimiters are not included in the trace. Therefore, session duration results which might be directly calculated from a tcpdump trace using TCP SYN and FIN markers, for example, had to be inferred here. In fact, we tried several different methods to infer session duration results, and concluded that while we cannot infer the session duration accurately enough, we can infer that the radio has a linear backoff registration scheme.

Figure 16 shows how often a given interval between successive registrations by a single radio occurs. The histogram shows exponentially decreasing peaks every 600 s, or 10 min. If the radio is not reset or if the radio does not change its primary poletop, the radio first registers 10 min after the initial registration, then 20 min, then 30, and so on. Rather than the typical exponential backoff found in networking protocols, a linear backoff is used.

Second, we wanted to investigate patterns of communication between radios. However, the trace has been anonymized to protect Metricom customers' privacy. As a result, we have no information about the other end of a radio's connection, or a mapping from radio to user, and can learn nothing about patterns of communications between radios.

Third, since we have a log from the nameserver only, we could not differentiate problems at the nameserver from prob-



Figure 16. Histogram of the number of times a time interval (between successive registrations by a radio) occurs. For example, an interval of length 600 s (10 min) occurs about 0.025% of the time, or about 150,000 times. Note the logarithmic scale for the *y*-axis, and that the tail of the graph, which extends to 4,444,118 s (74,068.6 min, 1234.5 hours, or 51.4 days) has been cut off to show the details.

lems in the network. For example, the trace includes three periods during which no radio registrations were logged: February 16, 6 a.m. through 1 p.m.; February 16, 5 p.m., through February 17, 4 a.m.; and February 17, 4 p.m., through February 18, noon. We do not know if these holes in the trace were due to a nameserver malfunction or a series of network outages. Metricom has since suggested that this outage is likely a failure in the logging mechanism rather than any network malfunction that would affect users.

Another example of such an ambiguity is a potential network reconfiguration on March 19. On that day, several poletops, previously unseen in the trace, first appear near a poletop at Metricom headquarters and then within the period of an hour, 5 p.m. to 6 p.m. on March 19, reappear at a new location in the Midwest. We do not know whether these poletops malfunctioned, or whether the network was in fact reconfigured. We suspect the poletop radios were being tested before installation in the Midwest.

Finally, the timestamps are the last problem we encountered in analyzing this trace. They are taken at the nameserver, which means that while clock skew is not a problem, network traversal time between the radio and the nameserver is not taken into account. Although we cannot determine the upper bound for this error, we can approximate the average error, since Metricom states that each packet will make no more than three wireless hops before reaching a wired access point. Traversing three wireless hops takes on average between 200 ms and 400 ms, which is less than the accuracy of the timestamp itself, making this a non-issue. These unanswered questions and ambiguities are all a direct result of not having the necessary data.

5.3. Clustering

The final lesson we learned while doing this analysis is the importance of parameter choice and distance function in both hierarchical and *k*-means clustering.

For example, when we use hierarchical agglomerative clustering to group poletops into locations, the choice of a cutoff distance is crucial. Hierarchical agglomerative clustering builds a tree by iteratively finding the two closest nodes without parents to become the children of a new parent node. Since the distance associated with each parent node is the distance between its two children, we can use this information to differentiate between different clusters. However, we found that there is no good static value to use as a cutoff between clusters.

We use a dynamic cutoff instead. Given all of the distances between children nodes, we look for the first distance which is greater than half a mile. We then look for an exponential increase (i.e., factor of two) over that distance. This new distance is the cutoff used to differentiate clusters of poletops into locations.

We then use *k*-means to refine the clusters. In *k*-means clustering, the number of clusters is chosen *a priori*, and the points are assigned to clusters iteratively. This algorithm is repeated for each radio.

Because we do not know *a priori* the number of locations a radio visits over the course of the trace, we cannot just use k-means or EM in this situation. Hierarchical agglomerative clustering helps us determine the number of locations for a radio. We also do not use EM in this situation, since the units in both parameters are the same.

The Metricom network already has a coarse granularity: no movement within a building can be detected. We make the network coarser by clustering poletops into locations, with the choice of a minimal distance affecting the resulting coarseness. Figure 17 shows how changing the minimal distance affects the number of locations a radio sees. For example, the most locations any radio visits using 0.5 miles as a minimal distance is 176. Using a minimal distance of 0.25 more than doubles this number to 428. We chose to underestimate the number of locations visited per radio, and thus underestimate the total mobility in the network.

An example of the importance of the choice of distance function in the clustering algorithm occurs in the second-stage clustering used to generate the results in table 3. Rather than using the standard Euclidean or Manhattan distance when clustering, we modified the standard Euclidean distance function for the clustering based on peak time usage such that the distance between midnight (0) and 11 p.m. (23) is 1 rather than 23. Finally, the two-stage clustering of radios to find patterns of mobility shows how parameter choice is crucial. When we clustered on poletops merely to find the logical grouping into locations, the choice of parameter was obvious: the latitudinal and longitudinal position of the poletop. For mobility, however, we had many choices in how to parameterize each radio.



Figure 17. Cumulative histogram showing how many radios visit a given number of total locations using a minimal distance of 0.25, 0.3, 0.4, or 0.5 miles. The tail of the graph, which extends to 428 locations, has been cut off to show the detail.

When we first clustered based on the overall mobility of the radio and on the frequency of radio usage, some other parameters we could have used were:

- The total number of registrations, since we might want to differentiate radios that are barely active during the trace. However, this parameter varied so much in value that it influenced the resulting clusters too much.
- The total number of poletops at which a radio registers. We felt that it would be better to use the location information derived from the poletops, since this is more reflective of the radio's actual movement than the poletops themselves.

We also could have chosen not to use all nine of the parameters we did use. Changing which parameters are used affects the results of the clustering.

Further, when we did the time-based second stage clustering and chose the overall peak times of day for the final dimensionality, we also could have chosen any of the following:

- 1. Peak time per location. Rather than two overall peaks, base the dimensionality on the total number of locations rather than the number of locations per active day. This choice does not work very well due to users who may visit many more total locations over the course of the trace than they visit per day, leading to unclear results. We could also look at peak time per day of the week, rather than overall peak times.
- 2. Durations rather than start times. Due to the linear backoff registration scheme, this parameter is hard to quantify.

6. Related work

We know of no other publicly available analysis of a metropolitan mobile network of this scale. However, there are studies of smaller mobile networks.

Researchers at CMU examined their large WaveLAN installation [4]. This study focuses on characterizing how the WaveLAN radio itself behaves, in terms of the error model and signal characteristics given various physical obstacles, rather than on analyzing user behavior in the network.

Another related set of papers is joint work from Berkeley and CMU [7]. The main paper in this set outlines a method for mobile system measurement and evaluation, based on trace modulation rather than on network simulation. This work differs from our own in several ways. First, the parameters they concentrate on deal with latency, bandwidth, and signal strength rather than radio movement. Second, their emphasis is on using these traces to analyze new mobile systems, rather than on understanding the current system. In this paper, our goals are to understand how people use an existing mobile system, with the eventual goal of providing parameters that could be used in simulating mobile networks in the future. Our current focus is on radio movement rather than radio behavior characteristics such as latency and bandwidth.

We previously performed a study of a combined wireless and wired network [5]. However, this study was limited in that only eight users participated rather than the 24,773 in our trace. Also, the Stanford study concentrated more on comparing which end-user applications were used in the wireless versus wired arena, and on determining the characteristics of the wireless network, such as latency and bandwidth.

7. Conclusion

Although the information in the traces limits the obtainable results and the results themselves are particular to the Metricom network, with its metropolitan-area coverage and high latency (compared to a local-area network), our analysis is a start on understanding how people take advantage of a mobile environment. We find that the more locations users visit on a daily basis, the closer together, on average, those locations are. In addition, the distance users move is a Gaussian distribution around the radius of the network. We also find that radios are used mostly during non-work hours, presumably due to users connecting to a faster network connection during work hours.

Given that many simulations of radio mobility assume a static time between radio movement or some other simple model of user movement, our results present a more realistic model of movement. We hope that these results can be used to guide simulations to yield projected network performance results based on observed movement.

8. Future work

The main body of future work needed is the comparison of this trace to other mobile network traces. It is possible, even probable, that at least some of the results are a product of the Metricom network and its features, such as its latency, network infrastructure, and trace information or lack thereof. For example, the latency in the Metricom network is high enough that people probably do not use it as they would an in-building wireless network such as WaveLAN. We would like to trace another network ourselves, and see how those results differ from these.

Acknowledgements

We thank Metricom, especially Mike Ritter and Loan Nguyen, for providing us with the traces. It is very unusual for a company to share its traces and we greatly appreciate it. Also, we thank Robert Bosch and Chris Stolte for their help with the graphs and the visualization, Lise Getoor and Uri Lerner for their help with the EM algorithm, Tamara Munzner for her help with LaTeX, and Petros Maniatis, Beth Seamans, Guido Appenzeller, Mema Roussopoulos, Ed Swierk, Lucas Pereira, and Andrew Beers for reading paper drafts. This work was supported in part by a generous gift from NTT Mobile Communications Network, Inc. (NTT DoCoMo), a grant from the Okawa Foundation, and a National Physical Science Consortium fellowship.

References

- [1] R. Bosch, C. Stolte, D. Tang, J. Gerth, M. Rosenblum and P. Hanrahan, Rivet: A flexible computer systems visualization environment, Computer Graphics 34(1) (February 2000) 68-73.
- [2] S. Cheshire and M. Baker, A wireless network in MosquitoNet, IEEE Micro (February 1996) 44-52.
- [3] A.J. Dempster, N.M. Laird and D.B. Rubin, Maximum likelihood from incomplete data via the EM algorithm, Journal of the Royal Statistical Society B 39 (1977) 1-38.
- [4] D. Eckardt and P. Steenkiste, Measurement and analysis of the error characteristics of an in-building wireless network, Computer Communication Review 26(4) (October 1996) 243-254.
- [5] K. Lai, M. Roussopoulos, D. Tang, X. Zhao and M. Baker, Experiences with a mobile testbed, in: Worldwide Computing and Its Applications,

Lecture Notes in Computer Science, Vol. 1368 (Springer, Berlin, 1998) pp. 222–237.

- [6] Metricom, Inc., http://www.metricom.com
- [7] B. Noble, M. Satyanarayanan, G. Nguyen and R. Katz, Trace-based mobile network emulation, Computer Communication Review 27(4) (October 1997) 51-61.
- [8] E. Rasmussen, in: Information Retrieval: Data Structures and Algorithms, eds. W.B. Frakes and R. Baeza-Yates (Prentice-Hall, Upper Saddle River, NJ, 1992) pp. 419-442.
- [9] S. Russell and P. Norvig, Artificial Intelligence: A Modern Approach (Prentice-Hall, Upper Saddle River, NJ, 1995).



Diane Tang is a Research Associate in the Department of Computer Science at Stanford University. Her interests include mobile networking, computer architecture, operating/distributed systems, and information visualization, especially with regards to level-of-detail issues. She is currently working on the Rivet and Polaris projects on interactive visual exploration of large datasets. Tang received an A.B. degree in computer science in 1995 from Harvard/Radcliffe University, and a Ph.D. in computer

science in 2000 from Stanford University. Tang is a recipient of the NPSC fellowship.

E-mail: dtang@cs.stanford.edu



Mary Baker is an Assistant Professor in the Departments of Computer Science and Electrical Engineering at Stanford University. Her interests include operating systems, distributed systems, and software fault tolerance. She is now leading the development of the MosquitoNet mobile and wireless computing project and the Mobile People Architecture. Dr. Baker received a BA degree in mathematics in 1984 from the University of California at Berkeley, and MS and PhD degrees in computer science in

1988 and 1994 also from U. C. Berkeley. Baker is a recipient of an Alfred P. Sloan Research Fellowship, a Terman Fellowship, an NSF Faculty Career Development Award, and an Okawa Foundation grant. She is a member of the ACM, the IEEE, and USENIX.

E-mail: mgbaker@cs.stanford.edu