

Capacity of Multi-service Cellular Networks

with Transmission-Rate Control: A Queuing Analysis

Eitan Altman
Projet MISTRAL, INRIA
Sophia Antipolis, France
altman@sophia.inria.fr

MOBICOM, ATLANTA, Sept. 2002

Outline

- Introduction
- Uplink static model
- Dynamic statistical model
- Overview of results
- Extending the model
- References

1 Introduction

There exist two main approaches to define capacities in wireless networks:

- **Pole capacity:** the maximum of the "load" that the system can handle (e.g. maximum number of mobiles). This is a **Static Notion**.

- **Erlang capacity:** defined with respect to a rejection limit ϵ . It is the "arrival rate" of connections that the system can handle while ensuring outages or drop of calls with probability less than ϵ .

- We call this the Erlang(ϵ) capacity.

Introduced in telephony networks in the beginning of the last century. Computed in CDMA context in A. M. Viterbi and A. J. Viterbi. Erlang capacity of a power controlled CDMA system. *IEEE Journal of Selected Areas in Communications*, pages 892–900, 1993.

- Consider here NRT traffic (e.g. data).
- NRT traffic does not have strict throughput requirements. Therefore the "pole capacity" (in terms of number of connections) is unbounded, if we can assign sufficiently small throughputs.
- What is the impact of smaller throughputs on Erlang capacity?
- If we decrease NRT rate, we can accept more calls but the transmission times are longer (in contrast to RT applications). Impact on Erlang capacity not clear.
- How to define a statistical type capacity when we can control the transmission rate?

Possible answers:

• Define the best Erlang capacity, where we optimize over different transmission rates.

(1st part of the talk)

• Erlang capacity is defined with respect to a given ϵ . Is there a version of the capacity that does not depend on ϵ ? (like Shannon capacity)?

(2nd part of the talk)

• Erlang capacity is related to REJECTION RATE.

If we can avoid rejections (by slowing transmission rates), one can define capacity related to OTHER QoS!!!

For example (ongoing work with Nidhi Hedge):

The arrival rates for which the expected delay is smaller than D

2 Uplink static Model

- Consider
- $K = \{1, \dots, k\}$ classes of service
- $M(s) =$ number of calls of type s , $\mathbf{M} = (M(1), \dots, M(k))$.
- $R(s) =$ transmission rate of class s
- The necessary received power at the BS for mobile of type s :

$$(1) \quad P(s) = \frac{N + I_{own} + I_{other} - P(s)}{\tilde{\Delta}(s)}, \quad s = 1, \dots, k.$$

where

$$\tilde{\Delta}(s) = \frac{E(s)WN_o}{R(s)}, \quad I_{own} = \sum_k^{j=1} M(j)P(j), \quad I_{other} = i \times I_{own}$$

(used when there is a constant number of mobiles).

• We can rewrite the power constraints as

$$(2) \quad P(s) = \frac{N + I_{own} + I_{other}}{\tilde{\Delta}(s)}, \text{ where } \Delta(s) = \frac{1 + \tilde{\Delta}(s)}{\tilde{\Delta}(s)}$$

$$\Leftrightarrow \tilde{\Delta}(s) = \frac{1 - \Delta(s)}{\Delta(s)}$$

$s = 1, \dots, k.$

• Recall: $I_{own} + I_{other} = (1 + i) \sum_{j=1}^k M(j)P(j).$

•The solution is

$$P(s) = \frac{N\Delta(s)}{1 - \sum_{j=1}^k (1 + \rho_j) M(j)\Delta(j)} \quad (3)$$

•The **Pole Capacity** is the polyhedron that makes the denominator vanish:

$$M^* = \{M : 1 = \sum_k (1 + \rho_j) M(j)\Delta(j)\}.$$

- A solution P is finite if and only if $M > m$ (in the Pareto sense) for some $m \in M^*$.
- A finite solution P exists if and only if all its components are finite. Hence when accepting a call when capacity is reached, not only its own QOS will not be respected, but also all ongoing calls will suffer.

Definition: Let \mathcal{M} be the subset of N^k for which $1 > \sum_{j=1}^k M(j)\Delta(j)$, and let

$$\eta = \max_{m \in \mathcal{M}} (1 + \eta) \sum_{j=1}^k m(j)\Delta(j). \quad (4)$$

We define the integer capacity M_B of the system as the boundary of \mathcal{M} for which adding a call results in an infinite power. It is the set of \mathbf{M} for which

$$\eta = (1 + \eta) \sum_{j=1}^k M(j)\Delta(j).$$

Definition: The blocking set M_B^j of class $j \in K$ is the subset of \mathcal{M} for which

another call of type j cannot be accepted:

$m \in M_B^j$ iff $m \in \mathcal{M}$ and $m + e_j \notin \mathcal{M}$, where e_j is the unit vector in direction j .

3 Dynamic statistical model

• Standard assumptions:

- Arrival of class s calls: Poisson process with parameter λ_s .
- Duration of calls of class s : exponentially distributed with parameter μ_s .
- Load of class s : $\rho(s) := \lambda_s / \mu_s$.

• The process of ongoing calls is a Markov chain with a stationary distribution $\pi_\rho(\mathbf{M})$.

• **Definition:** The Erlang capacity $EC(\epsilon)$ is the set of vectors $\rho = (\rho(1), \dots, \rho(k))$ s.t.

$$P_B(\rho) \leq \epsilon.$$

It is a set!

Theorem: The steady state probabilities of the Markov chain are

$$\pi_d(\mathbf{M}) = \frac{1}{G_d} \prod_{k=1}^s \frac{M^{(s)}_k}{d^{(s)}_{M^{(s)}_k}} i, \quad \mathbf{M} \in \mathcal{M}, \text{ where } G_d = \sum_{\mathbf{m} \in \mathcal{M}} \prod_{k=1}^s \frac{m^{(s)}_k}{d^{(s)}_{m^{(s)}_k}}. \quad (5)$$

The probability $P_s^B(d)$ that a class s class be blocked and the average blocking probability $P_B(d)$ are

$$P_s^B(d) = \sum_{m \in \mathcal{M}_s^B} \pi_d(m), \quad P_B(d) = \sum_k \lambda_s P_s^B.$$

Goal: Study the behavior of these probabilities and the Erlang capacity as a function of the transmission rate while keeping the same processes of arrival and size of files.

4 Overview of the Results

- Homogeneous case, single cell:
If we slow down the transmission rate by a factor a ,
- The blocking probability decreases,
- The Erlang capacity increases,
- As $a \rightarrow \infty$, the blocking probability tends to zero if $\rho > M_B$.

- Example:
- Consider a high speed application: $R = 160KB/sec$ (1.28 Mbps).
- Let $\Delta = 0.199$; the system is dimensioned to accept at most 5 simultaneous calls.
- The average quantity of transmitted information per call is 10KB (this is the measured average value over the Internet, see Sikdar et al, Perf. Eval, 2001).
- Hence the average call duration is $\mu^{-1} = \frac{10KB}{160KB/sec} = 62.4msec$, et $\mu = 16.03$.
- Assume that we wish to have a blocking probability of less than 1%. Then the Erlang Capacity is $p = 1.361$. We can accept $\lambda = p\mu = 21.8$ calls per sec.

| Slowing factor a | Δ | $\tilde{\Delta}$ | $M_{B,a}$ | Erl. Cap. EC(1%) | Arrival rate λ | call duration (msec) | Gain in % |
|--------------------|----------|------------------|-----------|------------------|------------------------|----------------------|-----------|
| 1 | 0.199 | 0.124 | 5 | 1.361 | 21.8 | 62.4 | 0 |
| 2 | 0.110 | 0.124 | 9 | 1.891 | 30.3 | 124.8 | 43.7 |
| 3 | 0.0764 | 0.0827 | 13 | 2.202 | 35.29 | 187.2 | 67.4 |
| 4 | 0.0583 | 0.0620 | 17 | 2.413 | 38.67 | 249.6 | 83.4 |
| 5 | 0.0473 | 0.0497 | 21 | 2.568 | 41.15 | 312 | 95 |
| 6 | 0.0398 | 0.0414 | 25 | 2.688 | 43.08 | 374 | 104 |
| 20 | 0.0122 | 0.0124 | 81 | 3.315 | 53.13 | 1248 | 144 |

Table 1: Gain in Erlang capacity by slowing transmission rates by a factor of a

- **Significance of results:** if we propose to subscribers services that differ in offered throughput, we can determine how to price per volume of traffic as a function of the *effective resources* that each service consumes.
- Example: We see that we double the capacity when dividing throughput by five.

- Hence a service five times faster should cost the double per traffic volume.
- Other factors that depend on the throughput may influence the pricing such as the equipment.

Single cell, heterogeneous case:

- Assume that the distribution of the size (number of packets) of class s call is exponentially distributed with parameter $\zeta(s)$
- Let $\nu(s) = \lambda(s)/\zeta(s)$,
- Let

$$\delta(s) = \frac{E(s)}{WN_o} \quad \text{and hence } \tilde{\Delta}(s) = R(s)\delta(s) \text{ and } \Delta(s) = \frac{R(s)\delta(s)}{1 + R(s)\delta(s)},$$

- If we slow down the throughput by a factor a , as $a \rightarrow \infty$, the blocking probability tends to zero if

$$\sum_k^{s=1} \nu(s)\delta(s) > 1.$$

- If $\sum_k^{s=1} \nu(s)\delta(s) > 1$ then for any \mathbf{R} ,

$$\sum_s P_s^B \delta(s)\nu(s) < \sum_s \delta(s)\nu(s) - 1 < 0.$$

Other results:

• A single cell with both RT and NRT traffic:

We compute the limit blocking probabilities and the capacity when slowing down

only the NRT traffic.

• The multiclass and multicell case

• We use a fixed-point approach to compute the blocking probabilities. It is based on replacing $I_{other} = i \times I_{own}$ by

$$I_{other} = iE[I_{own}].$$

• We compute the limit capacity as we slow down the transmission throughputs.

Fixed point approach

• We obtain

$$(9) \quad P(s) = \frac{1 - \sum_{j=1}^s M(j) \Delta(j) - Q}{N \Delta(s)}$$

where $Q = \sum_{j=1}^s E[M(j) \Delta(j)]$

• For each q (possibly different than 1), we obtain the probability

distribution of $M(s)$, $s = 1, \dots, k$ as before:

$$\pi^d(M) = \frac{1}{G^d} \prod_{k=1}^s \frac{G^d}{d^{M(s)}} \frac{i^{M(s)}}{M(s)!}, \quad M \in \mathcal{M}(q), \quad \text{where } G^d = \sum_{m \in \mathcal{M}(q)} \prod_{k=1}^s \frac{d^{m(s)}}{m(s)!}.$$

and where

$$\mathcal{M}(q) = \left\{ (m(1), \dots, m(k)) : \sum_{k=1}^j m(k) \Delta(j) > 1 - q \right\}$$

$$\mathcal{M}_s^B(q) = \left\{ (m(1), \dots, m(k)) \in \mathcal{M} : \sum_{k=1}^j m(k) \Delta(j) \geq 1 - b - \epsilon \right\}$$

• Define

$$F(q) = \sum_k \lambda_k E^q[M(j) \Delta(j)].$$

Then Q is the solution of the fixed point equation:

$$q = F(q).$$

(7)

- $F(q)$ is piecewise constant, and has thus discontinuities. Hence (7) need not have a solution. However, the set of values of q for which a solution to (7) does not exist has Lebesgue measure zero. A slight change in the value of q will yield a solution.
- $F(q)$ is nonincreasing in q which implies uniqueness of the solution to (7).

5 Extending the model

- Adding activity factors $\alpha(s)$ (between 0 and 1).
Approximation (Koo et al, 99): $I_{own} = \sum_{j=1}^k M(j)P(j)$ is replaced by $I_{own} = \sum_{j=1}^k \alpha(j)M(j)P(j)$.

- Imperfect power control:
multiply the $\Delta(j)$'s by some constants that depend on the standard deviation of the received signal to interference ratio.

- Exact analysis is not tractable, see [Viterbi+Viterbi, 93].

Dynamic throughput

- NRT share the bandwidth leftover from RT traffic.
- RT bandwidth can be reduced when the load is high.
E.g., UMTS will use the Adaptive Multi-Rate (AMR) codec that offers 8 transmission rates of voice between 4.75 kbps to 12.2 kbps.
- Ongoing work with Nidhi Hedge

6 References

- I. Koo, J.H. Ahn, J. A. Lee and K. Kim, "Analysis of Erlang capacity for the multimedia DS-CDMA systems", IEICE trans. Fundamentals, May 1999. (Includes Erlang capacity formula for multiservice systems using product form M/M/S/S type formulae).
- M. Meo, E. Viterbo, "Performance of wideband CDMA systems supporting multimedia traffic" IEEE Com Let, June 2001. (Includes Erlang capacity formula for multiservice systems using product form M/M/S/S type formulae).