

## Measuring a large network quickly

Ronald S. Burt <sup>\*,a</sup>, Don Ronchi <sup>b</sup>

<sup>a</sup> *Department of Sociology and Graduate School of Business, University of Chicago, Chicago, IL 60637, USA*

<sup>b</sup> *Don Ronchi Organizational Consultation, 12021 NW Fourth Street, Fort Lauderdale, FL 33325, USA*

We describe work in which we used three days of interviewing to identify and measure the network among 200 people significant in a complex production process. The capture–recapture strategy should be useful in other settings: (1) Conduct survey network interviews with people (informants) positioned in the study population such that their contacts overlap to provide recaptured relations. (2) Estimate reliability from data consistency across recaptures. (3) Triangulate relation response categories to assign quantitative scores to the categories. (4) Use reliability correlates to weight recaptured relations in the final network pooled across interviews. (5) Extrapolate from the known strengths of the captured relations to define the uncaptured relations.

Some networks take longer to measure than others. Dense networks reach closure faster, networks among easily available individuals can be measured faster, and small networks can be measured faster. There is no one size above which a network is deemed ‘large’ but we are here thinking of under 50 individuals as small and a network of hundreds or thousands as large. Measuring a network of more than a hundred individuals can require a month or more of rescheduled interviews or mail questionnaire reminders.

Time is not always available. For example, we were asked recently to measure the network structure of the production process in a manufacturing plant employing 2000 people. The network was known to extend beyond the plant to the headquarters of the firm owning the plant, but how it extended was unknown. Our work was a pilot intended to illustrate ways in which network analysis could be useful to help employees shorten the production cycle. We were expected to minimize disruption to employees and show value quickly. We spent a

\* Corresponding author.

day at the plant to learn the business and identify people to interview. We returned to conduct interviews, and were asked for summary results at the end of the first day. There were some disappointed to hear that summary conclusions would require analysis and more interviews. Patience ran out two weeks later, after a second day of interviews. Fieldwork was deemed complete. It was time to show value. We had 45 interviews. Interviews varied from 30 minutes to an hour, with about half of the time devoted to network data. We identified a population of 222 individuals significant in the production process and measured a network of connections among them. We used a capture–recapture strategy that should be useful in other settings; thus this paper.

## **1. Data collection**

The plant is part of the electrical industrial machines industry and is owned by one of America's largest manufacturing firms. Knowing the product is not essential for this paper, but it is important to note that the production process is complex and interactive. The plant's employees are distributed along and around what is in many ways a generic production process. Raw materials and some components are purchased from external vendors. These are combined at assembly points with components manufactured in the plant, leading to final assembly and product shipment.

Though in many ways generic, the production process is in other ways a study in extremes. Manufacturing tolerances are measured in millimeters, but upper limits on product size are set by the bore of railroad tunnels. Engineering and assembly are standardized, but idiosyncratic customer needs make each unit a 'customized' unit for which the details are worked out as the unit moves through the pipeline. It was clear on our first day that this was a production process of re-work and negotiation. Knowing the right people seemed essential to getting your work done.

### *Alternative data collection strategies*

Data collection could proceed in various ways. Indirect measures infer relations from data on joint involvements; relations between kinds of individuals involved in the same events (e.g. Burt and Minor 1983: Ch.

7, 8) or between specific individuals involved in the same events (e.g. Burt and Ronchi 1990). Indirect measurement is useful if aggregate features of a networks, such as a blockmodel or centralization, are to be traced over time or compared across social systems. It is less useful for a cross-sectional case study because you do not get data that describe the relation between two individuals. The data are their links with third parties, from which you infer their relationship with each other.

Similarly, we could not make use of 'network sampling' data collection strategies. These provide estimators of global properties such as the size or density of a network (e.g. Granovetter 1976; Frank 1978), or the typical person's contacts within the population network (Bernard *et al.* 1989; Freeman and Thompson 1989; Johnson *et al.* 1989), but there are few substantive research projects which have as their goal knowing that the network density of a study population is 0.44 or that the average person has 3000 contacts (Pool and Kochen 1978, discuss why this could be valuable). Our concern was to know who is connected to whom, and how well, in the production network of a specific organization (see Erickson 1978: 278–279, for heuristic illustration of the distinction between network structure and network sampling estimators).

The popular alternative for direct measurement is to collect sociometric choice data from a saturation or snowball sample of individuals in the study population (Coleman 1958; see Klovdahl 1989, for a broad review of sampling strategies). One strategy is to present a roster of the study population to each person, asking for the nature and frequency of their contact with each other person. This is impractical in a large network because people lose interest long before they finish searching through the list. A more traditional alternative strategy is to let respondents create their own lists of contacts in response to a name generator such as: "Who are your closest contacts?" Choices by respondent A define the strong connections in row A of the network. Strong connections in the other rows are defined by interviewing everyone in the study population.

Saturation strategies are not always practical, especially in a large network. As size increases, the boundaries of the network are likely to be ambiguous, it will take too long to interview everyone, and many of the interviews will be superfluous in the sense that they provide redundant information on the network's structure.

Snowball sampling is a popular alternative strategy. People in the study population define the boundaries of their network. Start with individuals in central positions. Ask for their key contacts. Interview the cited individuals, asking for their key contacts. Continue until some proportion of the individuals cited as key contacts have already been interviewed. There are virtues and drawbacks to snowball sampling (see Erickson 1978, for non-technical review; Goodman 1961 for an initial framing of technical issues; Snijders 1992 for contemporary review). For the purposes here, the self-defining boundaries of a snowball sample make the snowball strategy an attractive way to measure the core of a large network. The disadvantage under time pressure is that the strategy requires interviews with everyone in the snowball sample and intervals in the fieldwork during which you identify the next wave of people to be interviewed. We had to get in and out of the field quickly, with few opportunities to review our progress.

### *Capture–recapture*

We used saturation and snowball sampling, but replaced the usual sociometric questions with survey network questions. The primary difference involves eliciting data on the relations with and among people cited by the respondent (on the presumption that you will not be able to interview the people cited by the respondents, see Marsden 1990, for review discussion of survey network data). We began with name generators. The interviewed employee described the organization of his or her work – how the work is scheduled, where supplies come from, who depends on the output and soon – then named other employees involved in the work and most important to completing it. Name interpreter items followed. How often do you speak with him or her? How long have you known the person? Most important, respondents were asked to evaluate the strength of connections between people where connectivity was framed as follows:

We are interested in how people are connected along the production pipeline. People can be strongly connected in the sense that they often work together and work well together. At the other extreme, some people have a distant connection in the sense that they either rarely work together or wish that they didn't have to work together.

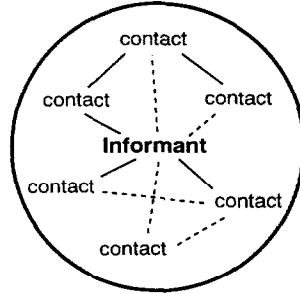
Respondents were asked to indicate their strongest connections on the list of people they had cited, and their most distant connections. Relations neither strong nor distant defined an intermediate, average, connection. After characterizing their relations with the listed people, respondents were asked to characterize relations between the listed people, strongly connected, distant, or somewhere in the middle. The result is a picture of the relations around the interviewed employee, a picture like the one at the top of Fig. 1, distinguishing strong connections (solid line), average (dashed line), and distant (no line).

These survey network interviews were conducted with people strategically placed along the production pipeline. There are five production units in the plant (shaded areas at the bottom of Fig. 1), each with a senior manager who reports to the plant manager. Each senior manager supervises two or more middle-managers. We focused on these middle-managers and the supervisors under them, the people who most closely direct production. We began with a saturation sample, every middle-manager and supervisor in A the five production units (Johnson's 1990 criteria 1 selection of informants). The snowball component was twofold: (1) During a break on the first day of interviewing, we cancelled scheduled interviews with a few people often cited by the individuals we had interviewed. We already had data on relations around often-cited people. Interviews with rarely cited production managers and supervisors were scheduled to replace the cancelled interviews. These were people around whom we did not have network data. (2) At the end of the day, we analyzed the network data to identify often-cited people in corporate headquarters and engineering who would be the target for our second wave of interviews.

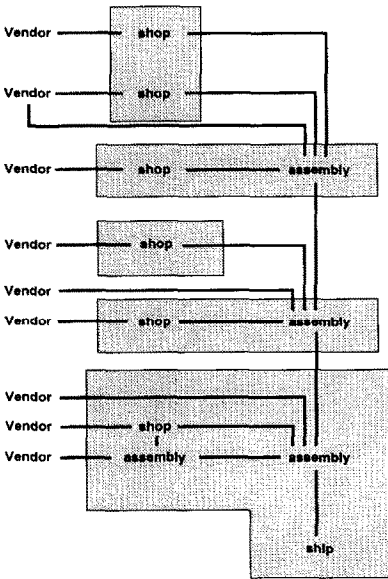
The data collection is an exercise in capture-recapture. The relation between two people is captured in an interview, then released for recapture in a later interview. Recapture marks relations central in the study population. Consistency across recaptures indicates data reliability. Individuals have two ways to appear in the network, one active, the other passive. The active role is to be one of the people interviewed. These are our informants. The passive role is to be cited, but not interviewed. These are contacts cited by one or more informants. For simplicity, we shall refer to these two roles from now on as informants and contacts.

We aimed for the diagram at the right of Fig. 1. The networks of

**Survey Network Interviews with Informants**



**Production Pipeline**



**Overlapping Networks Aggregate To Describe Social Structure along the Production Pipeline**

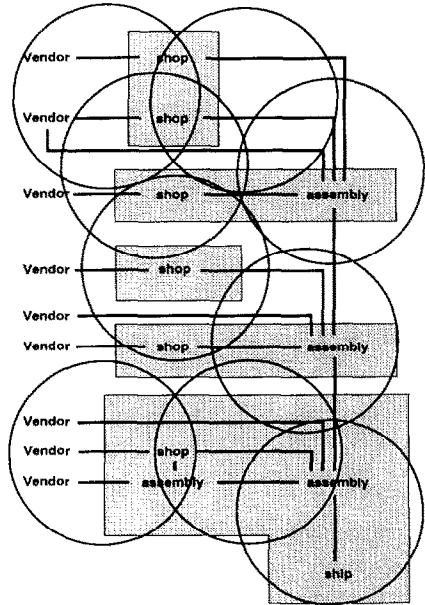


Fig. 1. Capture-recapture along the production pipeline.

people adjacent on the production pipeline overlap. Some contacts important to one person's work are important to the adjacent person's work. Sampling people for interviews is a balance between two concerns. To minimize cost and work disruption, informants should be far

enough apart in the production process to provide non-redundant network data. At the same time, they should be sufficiently close to provide multiple reports on strong relations so we get accurate measures of those relations and estimates of data reliability.

### *The production network*

The social structure of production is captured in a network of 222 people. We began with the managers at the top of the organization chart; the plant manager, senior managers who report to him, and middle-managers who report to the senior managers. Most of these core managers were cited as important contacts (49 of 59). The remaining ten are included to mark business functions disconnected from the production process. Around the 59 core managers, 132 other people in the plant or at corporate headquarters were cited at least once as an important contact, and another 13 were not cited as important but were cited by more than one person as a routine work contact. A further 38 people were never cited as important but were cited once as a routine contact. In 20 of these cases, the person was cited as a route to their supervisor, the cited person had a strong connection with their supervisor, and the supervisor was already among the 222 people in the network. These 20 citations were re-coded as citations to the supervisor. In six cases the person was not strongly connected to their supervisor. These six are included as individuals in the network. The remaining 12 people held supervisor rank and so are included as individuals in the network. In sum, the 222 person network comprised the 59 core managers in the plant, 132 other important contacts, 13 routine work contacts cited by more than one informant, and 18 other contacts.

Recapture frequencies are displayed in Table 1. The sum at the bottom of column one shows that 2121 relations were captured in our 45 interviews. Of these, 1621 were captured once (first row of table). We have only one informant's report on these. The second row shows that we have two reports on 344 relations. The last row of Table 1 shows the upper limit of nine reports on one relation. The 2121 relations were observed 2883 times. Observations are displayed in the three columns in the center of Table 1. The relationship captured nine times (bottom row of the table) was reported as average by one informant and reported as strong by the other eight informants.

Table 1  
2121 relationships, 2883 observations, and 1210 capture–recapture pairs

Number of relations	Times captured	Strength when captured			Loglinear z-score tendency to be strong	Capture–recapture pairs generated
		distant	average	strong		
1621	1	428	808	385	−6.58	0
344	2	142	354	192	−4.70	344
98	3	45	155	94	−2.82	294
29	4	10	72	34	−1.71	174
19	5	4	48	43	1.36	190
4	6	4	10	10	−0.62	60
4	7	4	10	14	0.20	84
1	8	0	2	6	1.39	28
1	9	0	1	8	2.02	36
2121		637	1460	786		1210

Strong relations tend to be recaptured. Loglinear parameters appear in the next to last column of Table 1. The parameters are z-scores that measure the tendency for relations to be reported as strong at each level of recapture. Relations captured only once tend not to be strong (−6.58 z-score,  $P < 0.001$ ). The tendency to be strong increases with recapture frequency.

The high recapture relations are not always obvious from an organization chart. Of the six relations captured seven or more times (bottom three rows of Table 1), none are distinctly prominent in the organization chart of the plant. Middle-managers under each of the five senior managers play different roles in the production process. Some are often cited as important to know. Others are rarely cited. The high-recapture relations connect each senior manager with his most active middle-manager, and connect some of those middle-managers with the plant manager.

A second point concerns further fieldwork. A snowball sample covers the study population when new citations lead back to people already interviewed. An indicator of our informant interviews covering the target population is the rate at which relations in new interviews are recaptures. Appropriately, relations captured only once tend to be distant connections in Table 1. Still, a quarter of them (385) are strong connections. This is a concrete indicator of what we knew at the



conclusion of our fieldwork; our data on some areas in the plant are lean. Strong connections in some segments of the production network are based on the report of a single informant.

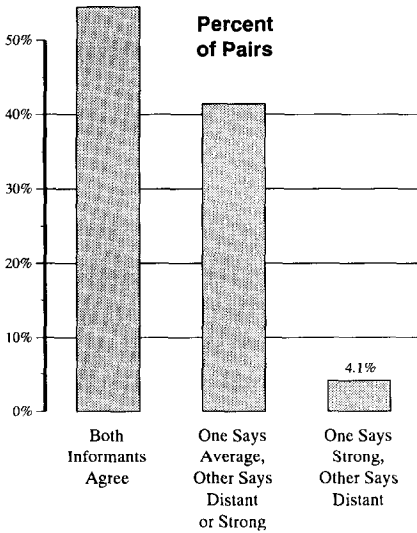
## 2. Reliability

With segments of the production network based on few data, data reliability is especially important. A virtue of the capture–recapture sampling is that reports can be compared across recaptured relationships. Reliability is high when multiple informants give similar descriptions of a relationship.

### *In the aggregate*

The last column in Table 1 lists 1210 paired observations for estimating reliability. We have only one person's description of the single-capture dyads, so there are no data for estimating reliability (0 in row one of the table). There are two descriptions of each of the 344 dyads captured twice. That yields, in the second row of Table 1, 344 paired reports to estimate reliability. When a relationship is captured three times, we have 3 pairs of reports (A versus B, A versus C, and B versus C), so the 98 relations captured three times yield, in the third row of Table 1, 294 paired reports for estimating reliability. And so on down the rows of the table.

The 1210 capture–recapture pairs are tabulated in Fig. 2 to illustrate two points about distinctions between strong, average, and distant relationships. First, descriptions are most often identical and rarely contradictory. The three cells on the main diagonal of the table are cases in which two informants give identical descriptions. These sum to 659, which is 54.5% of the 1210 pairs. Loglinear  $z$ -score effects for each cell of the table are displayed in Fig. 2. The only positive effects are in the main diagonal and they are strongly significant. With respect to contradictory reports, there are 50 instances of an informant saying that a relationship is distant when another informant says that it is strong ( $50 = 27 + 23$ ). Such contradictory reports are rare, they have strongly negative loglinear effects ( $-7.26$  and  $-7.14$ ), and a short bar in the graph (4.1%).



**Capture-Recapture Pairs**

	Distant	Average	Strong
Distant	61	85	27
Average	70	361	162
Strong	23	184	237

**Loglinear Z-Score Effects**

	Distant	Average	Strong
Distant	9.34	-1.85	-7.26
Average	-1.37	2.90	-0.86
Strong	-7.14	-0.42	10.23

Fig. 2. Reliability.

Second, average strength relations are the least reliable, but they are more easily distinguished from distant relations than from strong relations. The four cells adjacent to the main diagonal of the table are cases in which one informant said that a relation had average strength while the other reported it as distant or strong. These cases are a large proportion of the paired reports (41.4%), but the loglinear effects show that the only positive effect for average strength relations is the tendency for both informants to agree (2.90, versus the 9.34 reliability effect for distant and 10.23 for strong). The effects also show that average strength relations were less often confused with distant relations than with strong relations (effects of  $-1.85$  and  $-1.37$ , versus  $-0.86$  and  $-0.42$ ).

*Consistency across insiders and outsiders*

There is reason to expect higher reliability in some capture–recapture pairs than in others. Some are ‘informant–contact’ pairs in the sense that one or both of the people describing a relation is a participant in the relationship. In these pairs, your description of your relation with Joe is compared to an outsider’s description of the relationship. At least one informant is an insider. Other capture–recapture pairs are

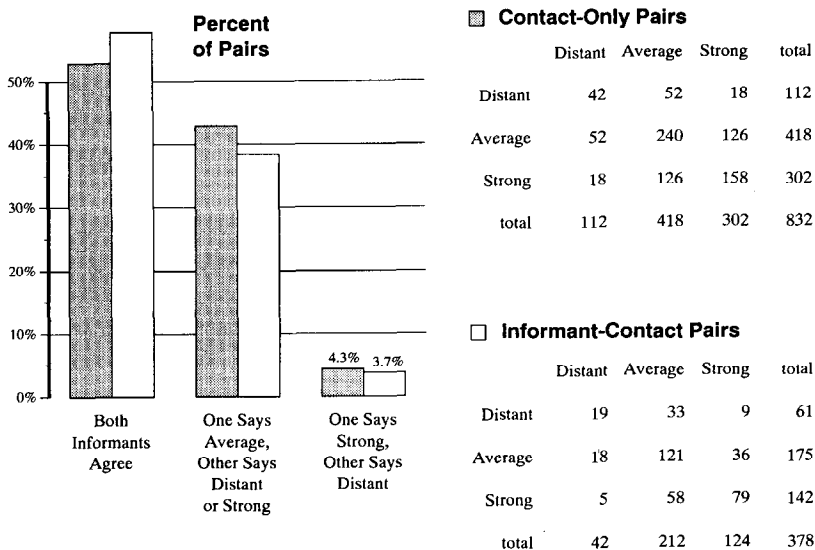


Fig. 3. Reliability when insiders and outsiders describe a relationship.

‘contact-only’ pairs in the sense that neither person describing the relationship is involved in it. Both informants are outsiders to the relationship. In a contact-only pair, one outsider’s description of your relation with Joe is compared with a second outsider’s description of your relation with Joe.

*Ceteris paribus*, insiders should be more reliable informants than outsiders. If this were not so, fieldwork would be much simpler. Instead of asking each person in a network to describe his or her relation with each other person, the network could be defined by asking one informant to describe the relation between each pair of people in the network.

Our third point on reliability, illustrated in Fig. 3, is that the data from outsiders are almost as reliable as the data from insiders. The figure contains a separate reliability tabulation for the 832 contact-only pairs and the 378 informant–contact pairs. There are a few cases where we captured both participants in the relationship; asking you about your relation with Joe, then asking Joe about his relation with you. With so few interviews in such a large population, however, these informant–informant pairs are rare and too few for reliable analysis,

so we have included the 25 informant–informant pairs in the tabulation of informant–contact pairs.

There are some expected differences. Informant–contact pairs are more likely to contain identical reports (57.7% vs. 52.9%), and less likely to contain contradictory reports (3.7% vs. 4.1%). Also, outsiders provide slightly cruder description. In the informant–contact tabulation, the informant’s description of his or her own relationship defines the row and the outsider’s description defines the column. The 18 in the second row and first column, for example, refers to 18 cases in which you said your relation with Joe was average strength but someone other than you or Joe said your relation with Joe was distant. Note the differences between the row and column sums. People describing their own relations often say that some are strong and others are distant, leaving slightly less than half for the residual category of ‘average’ strength (175 is 46% of the 378 relations). Outsiders are less likely to describe relations as strong or distant. They leave a higher proportion for the residual category (212 is 56% of the 378 relations).

These differences make substantive sense, but the important point is that they are negligible (12.19 chi-square with 8 d.f.,  $P = 0.14$ ). To distinguish strong, average, and distant connections, we get better, but not significantly better, data from people who are participants in the relations they are describing. This is an important conclusion because the capture–recapture data collection strategy depends on informants describing relations between other people.

### 3. Scaling connections

To pool recaptured relations across interviews and model the production network, we convert the response categories describing qualitative levels of connection strength into quantitative scores. Here again, the capture–recapture pairs can be useful.

#### *Response categories*

We have all relations characterized in terms of strong-average-distant, but we have richer data on the 440 relations between informants and contacts. The first column of Table 2 shows how relations are dis-

Table 2  
Paired reports for scaling the response categories <sup>a</sup>

Informant's view of relationship with contact	Frequency	Outsider's view of informant-contact relation				Informant's view of friend's relation with contact			
		distant	average	strong	total	distant	average	strong	total
Most distant	41	11 (3.7)	18 (-0.9)	5 (-3.1)	34	29 (5.0)	27 (-0.9)	10 (-3.4)	66
Distant	38	8 (3.2)	15 (-0.6)	4 (-2.7)	27	19 (3.9)	21 (-0.3)	7 (-3.0)	47
Average	199	18 (0.5)	121 (1.8)	36 (-2.3)	175	54 (0.9)	136 (1.1)	63 (-1.9)	253
Less strong	50	2 (-0.6)	24 (1.3)	9 (-0.3)	35	5 (-2.0)	24 (0.4)	23 (2.5)	52
Strong	46	3 (-0.7)	15 (-2.1)	31 (3.2)	49	5 (-2.1)	25 (0.5)	24 (2.6)	54
Strongest	66	0 (-2.1)	19 (0.5)	39 (3.5)	58	10 (-2.9)	42 (-0.6)	58 (4.6)	110
Total	440				378				582

<sup>a</sup> Loglinear z-score effects are given in parentheses and were computed for each tabulation separately. The most likely column response is boxed in each row. In some departments, there is a proprietary quality to an employee's strongest connections. The issue was identified by an unusually high frequency in cell (6,1) of the balance pairs tabulation. The 10 reported in the table is 10% of the row six average and strong relations, corresponding to 5 in the cell above being 10% of the row five average and strong ties. The observed frequency in cell (6,1) is actually 50, not 10. For the purposes of scaling, we put the disproportionate frequency aside, force the loglinear association model to ignore cell (6,1) when we scale response categories, and give the issue systematic attention in the next section.

tributed across six categories of connection strength. After an informant listed his or her contacts and the continuum from strong to distant was introduced (p. 94), the informant was asked to identify two people on the list, the informant's 'strongest' connection and 'most distant' connection. The informant was then asked if anyone else on the list was as close or as distant. This defined the two extreme categories of connection strength. The informant was asked if the remaining people on the list were all the same in distance between the two extremes. If yes, then the remaining people were assigned to the 'average' strength category. If no, which was usual, the informant assigned people to levels of strong and levels of distant connection. Any remaining people were assigned to the average strength category.

The six rows in Table 2 distinguish two categories of distant connection. Some informants distinguished up to four categories of

distant connection, but there are only eight relations in the third and fourth categories, and those eight have the same characteristics as the 30 relations in the second category (based on loglinear models of tabulations with frequency and other variables). The 38 relations in all three response categories of less distant connection are combined in a 'distant' category.

Three categories of strong connection are distinguished in Table 2. Some informants distinguished up to seven categories of strong connection. Again, there are few relations in the four categories beyond the third (19 in total), and they have the characteristics of the 31 relations in the third category of strong connection. The 50 relations in the five categories of weakest strong connection are combined in a 'less strong' category.

#### *Two associations between response categories*

Two associations between response categories are displayed in the table. The first panel is a tabulation of capture–recapture pairs. These are the 378 informant–contact pairs in Fig. 3, now in the middle of Table 2. The rows of Table 2 describe an informant's report of his or her relation with a contact and the columns describe the strength of the relationship as reported in a separate interview with an informant outside the relationship. Rows and columns are paired measures of the same relationship, so they can be used to scale one another.

The second panel in Table 2 is a tabulation of relations paired within interviews. Refer to the sixth-row contacts as 'friends.' These are the informant's closest contacts. Columns in the second Table 2 panel describe the informant's view of his or her friend's relationship with the contact cited for the rows. Under the principle of structural balance, the informant's relation with a contact should be the same as his or her friend's relation with the contact, friends of my friends are my friends, enemies of my friends are my enemies. To the extent that balance exists in the relationships, the rows and columns are paired measures of the same relationship and can be used to scale one another. Balance pairs are a useful basis for scaling survey network data, which typically do not provide capture–recapture pairs (Burt and Guilarte 1986; cf. Faust and Wasserman 1993, on scaling two or more kinds of relations under a multiplexity presumption).

## Results

The two associations require different assumptions about the response categories. To scale across the capture–recapture pairs, we assume that different informants refer to the same underlying continuum of strong to distant when they report on a relationship. To scale across the balance pairs, we do not have to make that assumption because both reports are from the same interview. Reports from different informants can refer to different understandings of the continuum from strong to distant connection. However, scaling across the balance pairs assumes that relations are balanced and that need not be true (discussed in the next section).

Given the differences between the two kinds of pairs, it is reassuring to see their similar results. First, the distribution of informant–contact relations in each tabulation represents the distribution of informant–contact relations in the total sample. There are no significant differences between the three columns of row marginals (9.93 chi-square, 10 d.f.,  $P = 0.45$ ). Second, we used loglinear models to identify the most likely column response in each row and reported the loglinear z-score effects in parentheses. The first two rows show that when an informant reports a relationship as ‘distant,’ the paired report tends to be distant and tends not to be strong. Further, the most distant relations are more likely to be corroborated as distant (loglinear effects of 3.7 versus 3.2 in the first and second rows of capture–recapture pairs, and 5.0 vs. 3.9 in the balance pairs).

The third row of Table 2 shows that when an informant reports a relationship as ‘average’ in strength, the most likely paired report is average. The reliability analysis showed that this category is more ambiguous than strong or distant. The point is here in a different form. Neither tendency for average to be paired with average is significant.

The bottom three rows of Table 2 show what happens when an informant reports a relation as strong. The extreme categories of ‘strong’ and ‘strongest’ are significantly often paired with reports of strong. The strongest relations are the more likely to be corroborated as strong (loglinear effects of 3.5 versus 3.2 in the sixth and fifth rows of capture–recapture pairs, and 4.6 versus 2.6 in the balance pairs). ‘Less strong’ relations are somewhere between average and strong. In the capture–recapture pairs, an informant’s less strong relations are

most likely to be viewed as average. In the balance pairs, less strong relations are significantly often paired with strong relations.

To scale the response categories, we fit a one-dimensional loglinear association model to the pooled frequencies (Goodman 1984). The model generates the loadings on the left in the below list from which we obtained the network scores on the right; scores for the column response categories:

distant	- 0.733	0.00
average	0.055	0.56
strong	0.678	1.00

and for row responses:

most distant	- 0.506	0.17
distant	- 0.465	0.20
average	- 0.137	0.45
less strong	0.143	0.66
strong	0.373	0.84
strongest	0.591	1.00

where the network score for response category *i* is (category *i* loading +0.733) divided by (maximum loading +0.733), the 0.733 is the minimum loading for any response category, and maximum loading is the row maximum for row loadings or the column maximum for column loadings. Network scores vary between the convenient extremes of 0 for the most distant connection, and 1 for the strongest connection. Figure 4 displays quantitative distances between the quali-

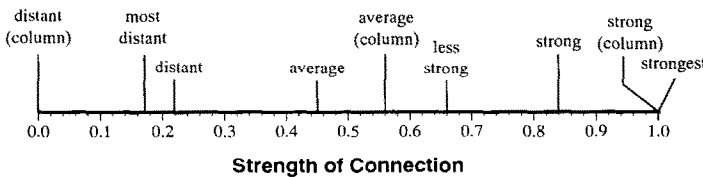


Fig. 4. Response categories positioned on the continuum from distant to strong connection.



tative response categories. Distant relations cluster at the bottom of the scale. Less strong relations cluster with the average strength relations. Stronger relations cluster at the top of the scale. We are not proposing that these distances generalize to other study populations. Maybe they do; more likely, they do not. The task here is to key the qualitative response categories, which do apply across populations, to a quantitative strength of connection that the categories represent in our specific study population.

### *Tertius bias*

The tabulation of balance pairs was distorted by an effect that can be termed *tertius* bias. Detecting the bias and tracing it back to certain functions within the organization had important substantive implications, but we mention it here because we had to adjust for it to fit the association model to the tabulation of balance pairs (right panel in Table 2). The *tertius gaudens* is a social structural definition of an entrepreneur and appears in several forms in network analysis (Burt 1992). The *tertius gaudens* is the 'third who benefits' by brokering contact between other people. *Tertius* bias refers to people describing social structure as if they were the *tertius*, exaggerating the extent to which they brokered the connection between their closest contacts.

The bias is illustrated in Fig. 5. The bottom graph is taken from the tabulation of balance pairs in Table 2. Given informants and their close contact 'friends,' when an informant has a most distant connection to a contact, 44% of the friends have a distant connection to the contact (29 divided by 66 in row one of Table 2). The stronger the connection between informant and contact, the less likely that friends are distant from the contact. When an informant has a strong connection with a contact, only 9% of the friends have a distant connection with the contact (5 divided by 54 in row five of Table 2). The graph at the top of Fig. 5 is taken from a tabulation of balance pairs in another study population (Burt 1992: Ch. 4). In both graphs, the grey bars, indicating distant connection between contacts, decrease as the informant's connection with both contacts strengthens.

*Tertius* bias is indicated by the black bar over 'strongest' connections. The grey bar over strongest connections is based on the frequency of 10 in cell (6, 1) of the balance pairs tabulation in Table 2.

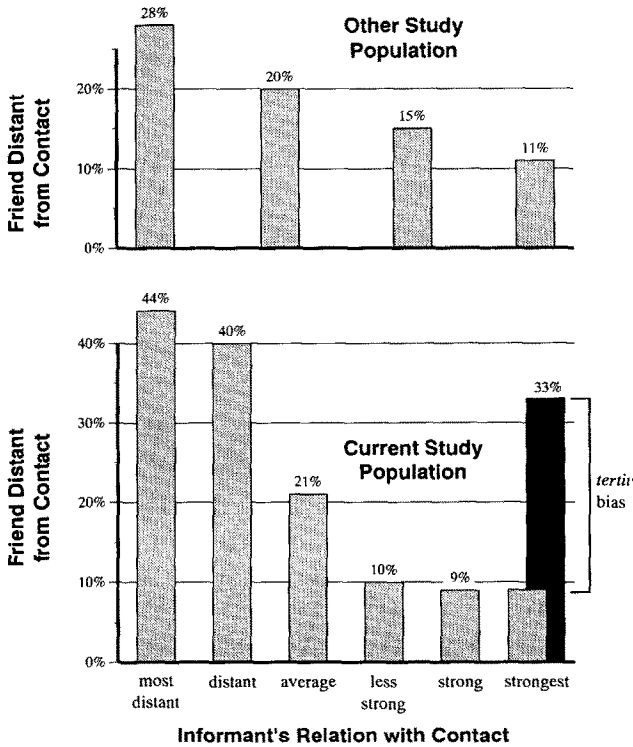


Fig. 5. *Tertius* bias.

This is an imputed frequency (see note (a) to Table 2). The actual frequency is 50, which defines the black bar in Fig. 5. Fifty is an unexpectedly high frequency of distant connections between people who should be close. There is little in the Table 2 tabulation of balance pairs that is not described by the continuum in Fig. 4 between distant and strong connections between employees (2.68 chi-square, 4 d.f.,  $P = 0.61$ ). However, the single continuum assumption is clearly rejected if the 10 frequency is changed to its true value of 50 (26.10 chi-square, 4 d.f.,  $P < 0.001$ ). A second dimension is required just to describe the high cell (6, 1) frequency of distant connections between people who should be close. Across relations, in other words, the strength of connection between contacts increases with the strength of an informant's connection to the two contacts. But the trend stops and reverses for the informant's closest contacts in this study population.

The black bar shows that an informant's strongest connections tend to be reported as distant from one another. This indicates *tertius* bias. There is a proprietary quality to the informant's role. His or her role in the network is to broker the connection between the contacts.

Why not believe the informant? Perhaps he or she is a *tertius*, necessary to make the connection between certain other people. We see the proprietary connections in the black bar as response bias because they are inconsistent with broader patterns in the network data. They are (a) inconsistent with balance theory evidence here and elsewhere of stronger connections between a person's closer contacts, (b) unnecessary in that people can have the entrepreneurial opportunities of disconnected contacts without the disconnections being concentrated among their closest contacts, and (c) uncorrelated with the actual volume of disconnections between employee contacts.

The data reported at the top of Fig. 5 were collected, with a survey instrument similar to the one used here, in a large American computer firm. Fewer distinctions were recorded between levels of connection between informant and contact (four bars rather than six) and the informants occupied positions higher in their organization, but it was a similar exercise of sorting relations into strong, average, and distant. A person's strongest connections in the other study population tend to be connected with each other (11% distant). There is no evidence of *tertius* bias; the final grey bar is not disproportionately high. Nevertheless, there are many disconnected contacts and strong network effects can be traced to the disconnections (Burt 1992: Ch. 4).

Further, the black bar disconnections are uncorrelated with disconnections more generally. Using the final network of pooled relations (described in Section 4 below), we measured the extent to which each of the 222 people in the production network is connected with people disconnected from one another. We measured *tertius* bias as the number of times that an informant reported distant connection between two of his or her closest contacts. Correlations are in the first column below. The second column lists correlations with a dummy variable; 1 if an informant showed any *tertius* bias (reported any distant connections between his or her closest contacts), 0 otherwise:

0.01	0.07	number of contacts
-0.01	-0.09	number of non-redundant contacts

- 0.05 – 0.05 density among informant's contacts
- 0.07 – 0.07 network constraint

The correlations are negligible. *Tertius* bias is uncorrelated with (1) number of contacts, (2) size adjusted for disconnections among contacts, (3) density, and (4) the constraint of coordinated contacts (see Burt 1992: Ch. 2, for review of the non-redundancy and constraint measures).<sup>1</sup>

Our sense of the proprietary connections as a *tertius* bias is further reinforced by where they occur. The inference from Fig. 5 is that all employees are to some extent proprietary about their closest contacts, brokering the connection between two people who would otherwise be distant from one another. However, the tendency varies by function. Employees in engineering and production control are significantly more prone to *tertius* bias.

Consider Fig. 6. Four broad functions are distinguished; production shops are the five shaded areas in Fig. 1 where the product is manufactured, production control refers to the people who control materials and quality, production support refers to people who facilitate production (human resources, facilities and equipment maintenance, community and environmental issues, and so on), and engineering refers to the people who design and draft specifications for the product. Of the 222 people in the production network, all but 32 are in one of these four functions. The remaining 32 are contacts in corporate headquarters.

The white bars in Fig. 6 show the proportion of people in each function who have entrepreneurial networks. This is a crude distinction between clique employees who have a few contacts and those few are strongly connected with each other, versus entrepreneur employ-

<sup>1</sup> We also checked capture–recapture reliability. *Tertius* bias involves two kinds of relations, an informant reporting maximum strength connections with two contacts and a minimum strength connection between the contacts. Reports from other informants could be used to corroborate these reports, but the relations were rarely recaptured (8 of the informant–contact relations and 1 of the distant connections between contacts). Because of the small numbers involved, we resist the urge to interpret the failure to recapture. The *tertius* bias is concentrated in engineering, and our data on engineering, collected on the last day of fieldwork, are not as rich as our data on production.

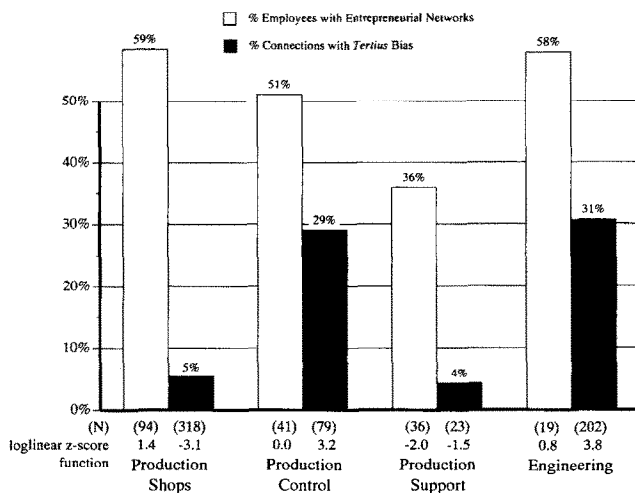


Fig. 6. Locating entrepreneurs and *tertius* bias.

ees who have many, disconnected contacts. The employees with entrepreneurial networks are defined as anyone with a below-median network constraint score (see Burt 1992: Ch. 2). The point illustrated is that there are entrepreneurial employees in all functions. The tendency for an employee to have an entrepreneurial network of disconnected contacts is independent of the function in which he or she is employed (5.52 chi-square, 3 d.f.,  $P = 0.14$ ).

The same is not true of employee tendencies to describe relations as if they had entrepreneurial networks. *Tertius* bias is indicated by the black bars in Fig. 6. The black bars indicate the proportion of balance pair connections through each function that show *tertius* bias.<sup>2</sup> The hypothesis of *tertius* bias being independent of function is clearly rejected (78.32 chi-square, 3 d.f.,  $P < 0.001$ ). The loglinear z-scores at the base of the bars in Fig. 6 show that *tertius* bias is

<sup>2</sup> There are only 50 balance pair connections that show *tertius* bias in the extreme sense of distant connection between an informant's closest contacts (black bar in Fig. 5). To improve comparisons across functions, we expanded the category to include four cells in Table 2 (the three cells of distant connection between an informant's strong contacts, and the category of average connection between an informant's strongest contacts). Of the 622 balance pair connections in Table 2 (582 in the table plus 40 excluded as explained in table 2 note (a)), 102 show *tertius* bias in this expanded sense.

concentrated in production control (3.2 z-score) and engineering (3.8 z-score).

In other words, there are employees in each function who have entrepreneurial networks in the sense that they broker connections between other employees, but it is in production control and especially engineering (larger z-score, much larger number of affected connections) that employees feel proprietary about the connections they broker. This is tied to the way work is carried out in the plant. In production, employees broker connections to get on with their work. In production control and much of engineering, brokering connections is their work. The tendency for employees in production control and engineering to hold the view that they were the only viable connection between their closest contacts had special significance for us because these two functions were most often blamed by other employees for production delays and re-work. That too encouraged us to put aside the *tertius* bias as a response effect when we scaled levels of connection strength to pool relations across interviews.

#### **4. Pooling across interviews**

One in four captured relations was captured more than once. Most of the 500 recaptured relations in Table 1 were captured two or three times, but some were captured several times. We have to pool the recaptured relations so that we have one relationship between each pair of employees in the network. There are options. To emphasize disconnections in the production network, the pooled relation between two employees could be set equal to the weakest reported relationship between them. Or, setting the pooled relation equal to the strongest reported relation would minimize disconnections in the network.

We had no substantive reason to emphasize or minimize disconnections, so we averaged relations across recaptures. If one person said a relation was strong (1.00) and a second person said it was of average strength (0.56), the unweighted average across the two reports is 0.78.

But we do not have equal confidence in every report. An informant who knows two people well is more likely to give an accurate report than someone only vaguely acquainted with the two people. To reduce error around the pooled relations, we weighted reports before averag-

ing them. This required other data on the relations and knowing how reliability is correlated with the other data.<sup>3</sup>

### *Other data on the relations*

We gathered data on how often informants spoke with the people they cited. If we know how the strength of a relationship is naturally associated in the population with specific levels of contact frequency, we can better define how often people should be asked to check in with one another if they are to be strongly connected in a reorganization. Frequent contact contributes to connection strength as we defined it for respondents (p. 94), but what it means to have frequent contact and how strength varies across specific levels of frequency are unknown.

Relations between managers often survive on a weekly or monthly rhythm in the sense that you tend to meet close contacts about once a week (or once a month) in one context or another; this week in a bi-weekly committee meeting, next week in the hall outside your office. Relations in this study population survive on a rhythm of daily encounters. The pie chart at the top of Fig. 7 shows that over half of the contacts are met once a day (55%). A quarter are met at least once a week. Among the 440 relations are 232 that the interviewed people deemed especially important. Two-thirds of these key contacts are met every day (64%).

The bar graph at the bottom of Fig. 7 shows how frequency is associated with strength. Strength is clearly contingent on frequency. The likelihood ratio chi-square of 125.7 with 10 degrees of freedom for independent frequency and strength can be rejected at beyond the 0.001 level of confidence. The association is concentrated in the

<sup>3</sup> We thought that regression to the mean would be a problem. It is unlikely that an extremely strong or extremely distant connection will be consistently strong or distant across recaptures. Now and again someone will report either extreme as an average strength relationship. Averaging across recaptures should make it difficult for recaptured relations to reach extreme values. The weighted averages incorporating data on frequency and duration would help preserve extreme values, however, regression to the mean turned out not to be a problem. In Table 1, 17% of the recaptures are 'distant,' 52% are 'average' and 32% are 'strong' relations (column sums excluding row one). We computed unweighted averages across recaptures. These 500 pooled relations in the second through the ninth rows of Table 1 are 19% 'distant,' 47% 'average' and 34 'strong' (where average is a relation greater than 0.3 and less than 0.7; cf. Fig. 4). Extreme categories are as present in the pooled data as in the original reports. In other words, reliability is sufficient to preserve extreme categories across recaptures.

**Cited contacts tend to be met every day. Those met less than weekly tend to be viewed as distant connections.**

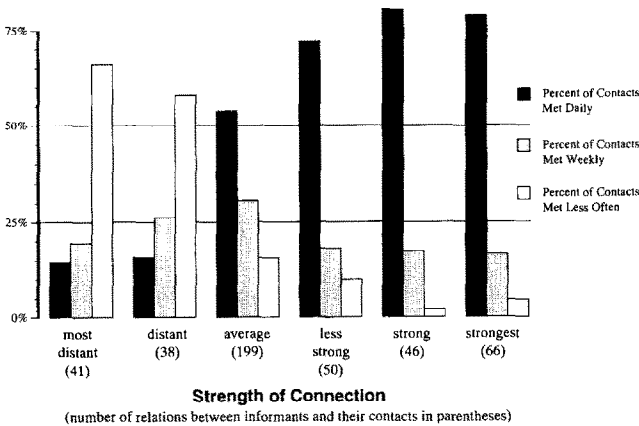
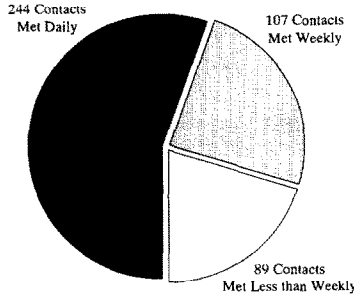


Fig. 7. Frequency and connection.

distinction between distant and other connections. A large proportion of distant contacts are met less than weekly (two left-most white bars, loglinear z-score effects of 6.1 and 5.1), but very few average or stronger contacts are met less than weekly. A small proportion of distant contacts are met every day (two left-most dark bars, loglinear z-score effects of -4.5), but the majority of average and stronger contacts are met daily.

Strong relations involve daily contact, and contacts met less than weekly tend to be viewed as distant, but long acquaintance can substitute for frequency. Many relations in the plant go back for 20 years. A strong connection between people who have known each other for many years need not depend on frequent meetings. The question is whether strong connections exist between recent acquaintances. If strong connections depend on long acquaintance and pro-



Table 3  
Mean years of acquaintance

	Daily	Weekly	Less often	Total
Most distant	6.5	2.9	2.3	3.1
Distant	3.3	9.4	2.4	4.4
Average	5.8	4.2	4.6	5.1
Less strong	4.1	4.8	2.0	4.0
Strong	4.8	3.4	1.9	4.5
Strongest	7.7	4.0	1.0	6.8
Total	5.7	4.5	3.3	4.9

duction process depends on strong connections, then changes to the production process depend on employees who have known each other for a long time.

The employees seem more adaptive than that. On average, the people in the cited relations have known one another for 4.9 years. The average varies with how often they meet each other and connection strength, as displayed in Table 3. Analysis of variance around these means shows that years of acquaintance do not vary systematically with connection strength but do vary significantly with frequency ( $F$ -ratio of 1.7 with 5 and 433 d.f. for no difference across the strength categories,  $P = 0.14$ , versus  $F$ -ratio of 7.1 with 1 and 422 d.f. for the hypothesis of no difference across frequency levels,  $P = 0.008$ ). A closer look at the data with loglinear models reveals a significant change in the first year. There are no significant tendencies for a relation to be strong if it is between people who have known one another for three, five, ten, or 20 years, but relations less than a year old tend to be viewed as distant (loglinear  $z$ -score effect of 2.6 for most distant, 2.0 for distant).

In sum, strong connections in the plant are rare between people who meet less than once a week or who have known one another for less than a year. A year acquaintance establishes a relationship, after which connections develop as strong or distant regardless of the years for which two people are acquainted. Strong relations tend to involve daily contact.

#### *Effect on reliability*

We have three indicators of an informant's familiarity with the relationship between two contacts: (1) The stronger the relations with

each contact, the more likely that the informant knows contact feelings for one another. (2) The more frequent the meetings with each contact, the more likely the informant has witnessed interaction between the contacts. (3) An informant with established relations to both contacts is more likely to know their relationship than someone who only met them recently. *Ceteris paribus*, we expect more reliable reports from informants with strong-frequent-established relations to both contacts (cf. Romney and Weller 1984: 66–69; Freeman *et al.* 1989). The empirical question is how much, and in what way, reliability depends on these factors.

Table 4 contains the answer. The 1210 capture–recapture pairs of reports are tabulated, informant–contact pairs separate from contact-only pairs. The two columns of ‘identical’ reports are instances of two informants both describing a relationship as distant, average, or strong. These are the diagonal elements of the tables in Fig. 3. Of the 378 informant–contact pairs, 219 are identical reports. Of the 832 contact-only pairs, 440 are identical reports. Any difference between the two reports moves the pair into the ‘Different’ columns in Table 4.<sup>4</sup>

Rows of Table 4 distinguish levels of informant familiarity with the people whose relationship is described. The first row refers to capture–recapture pairs in which both informants have strong-frequent-established relations with the two people connected by the described

<sup>4</sup> It might seem odd, having scaled the response categories, that we return to a reliability dichotomy between identical and different. We do so because quantitative criteria give us the same results as the qualitative dichotomy in Table 4 and the reliability correlates have discrete rather than continuous effects (footnote 5). The scaled relations have a standard deviation of 0.3361 across the 2121 captured relationships. For a quantitative measure of the difference between the reports in a capture–recapture pair, we divided the absolute difference between the two reports by the 0.34 standard deviation. If one informant said a relation was strong (1.00) and the other said it was less strong (0.66), the absolute difference between their reports is 0.36, which is a 1.06 z-score difference when divided by the 0.34 standard deviation. The average z-score difference across the 1210 capture–recapture pairs is 0.75 with a 0.78 standard deviation. We studied how this criterion, and others based on the magnitude of differences, varied across levels of strong-frequent-established relations (later simplified to the eight levels in Table 4). The most systematic covariation occurs at the bottom of the z-score scale, through the interval separating identical reports from not identical reports. That is the difference measured by the dichotomy in Table 4 between identical and different.

relationship. The bottom row refers to pairs in which both informants have distant-infrequent-recent relations with the two people.<sup>5</sup>

Two points are illustrated in Table 4: First, there is no significant insider–outsider difference. *Ceteris paribus*, reliability should be higher in the informant–contact pairs because one of the informants is describing his own relationship. We used Fig. 3 to show that reliability is higher in these pairs, but not significantly. Table 4 adds the further result that both kinds of reports are similarly affected by the structural conditions that make an informant competent. The structure of the informant–contact tabulation in Table 4 is not significantly different from the structure of the contact-only tabulation (5.68 chi-square, 8 d.f.,  $P = 0.68$ ).

In other words, strong-frequent-established relations improve reliability in both kinds of pairs in the same way. Reliability is keyed to an

<sup>5</sup> We began with detailed levels rather than the dichotomies in Table 4. The two informants in a capture–recapture pair each have a relation with the two people whose relationship is being described. We measured strong-frequent-established for each of the four informant relations as follows: strength is the network score for strength of connection with the informant (Fig. 4); frequency is 3 for daily, 2 for weekly, 1 for less than weekly; established is 1 if the informant has a year or less acquaintance with the contact, 2 if more than a year acquaintance. In informant–contact pairs, we coded the informant’s relation to him or herself as maximum strength. Multiply the four variables for the four informant relations in a capture–recapture pair to create a product variable: ESTABLISHED – The product of the established relation variables varies from 1 ( $1 \times 1 \times 1 \times 1$ , if both informants are a year or less acquainted with each contact) up to 16 ( $2 \times 2 \times 2 \times 2$ , if both informants have established relations with both contacts). Across levels, we found that reliability was high at level 16, but did not increase systematically up the lower levels. Capture–recapture pairs are sorted in Table 4 into two kinds under established; those in which both informants have established relations with both contacts (‘yes’), vs. those in which either informant is a recent acquaintance to either contact. FREQUENT – The product of the frequency variables varies from 1 ( $1 \times 1 \times 1 \times 1$ , if both informants speak less than weekly with both contacts), up to 81 ( $3 \times 3 \times 3 \times 3$ , if both informants speak daily with both contacts). Reliability is high in the three highest levels of frequency: daily-daily-daily-daily, daily-daily-daily-weekly, daily-daily-weekly-weekly. If either informant speaks weekly with both contacts, or less than weekly with either contact, reliability decreases. The decrease is not systematic down the product variable. Reliability rises and falls between adjacent levels of less frequent meeting. Capture–recapture pairs are sorted in Table 4 into two kinds under frequency; the above three high frequency conditions (‘yes’), vs. lower frequencies. STRONG – The product of the strength of connection variables varies from 0.0008 ( $0.17 \times 0.17 \times 0.17 \times 0.17$ , if both informants have a ‘most distant’ connection to both contacts), up to 1 ( $1 \times 1 \times 1 \times 1$ , if both informants have ‘strongest’ connections with both contacts). Reliability varies substantially at all levels of this product variable (there is only a 0.07 between reliability and the product variable), but it tends to be higher when three of the four informant relations are some level of strong connection. The capture–recapture pairs are sorted in Table 4 into two kinds under strong; those where at least three of the informant relations are strong (‘yes’), versus all lower strengths of connection.

Table 4  
Reliability correlates

Relations from informants to reported relationship			Loglinear z-score for identical reports	Capture-recapture paired reports				total
strong	frequent	established		informant-contact		contact-only		
				identical	different	identical	different	
yes	yes	yes	3.07	95	38	31	14	178
yes	yes	no	0.73	19	10	18	13	60
yes	no	yes	0.47	4	3	6	3	16
yes	no	no	-0.01	6	4	14	11	35
no	yes	yes	-0.78	50	38	106	93	287
no	yes	no	-0.04	13	11	60	44	128
no	no	yes	-2.09	11	16	41	45	113
no	no	no	-2.59	21	39	164	169	393
			2.89	219	159	440	392	1210

informant having information on a relationship, not to the informant being personally involved in the relationship. Personal involvement is one way to obtain information on a relationship, and reliability is slightly higher in informant-contact dyads. The more significant factor is how often an informant has occasion to observe the relationship he or she is describing. Just as an informant can more reliably describe his or her own strong-frequent-established relations, the informant can more reliably describe relations between contacts with whom he or she has strong-frequent-established relations. This further assuages our concern about reliance on informants in the capture-recapture data collection. Variation in reliability can be traced to variation in relations, but not to the difference between insiders and outsiders describing the relations.

Our second point from Table 4 is the manner in which reliability changes across capture-recapture pairs. The column of loglinear z-scores in the middle of Table 4 describes the tendency for identical reports to occur in the conditions described by each row. In the aggregate, reports in the capture-recapture pairs tend to be identical across conditions (2.89 z-score). The aggregate tendency is significantly higher when both informants have strong-frequent-established relations with the contacts between whom they are describing a relationship (3.07 z-score). At the other extreme, the tendency for reports to be identical drops significantly if the informants have

distant, infrequent relations with the contacts. These are the bottom two rows in Table 4. Between the two extremes, reliability does not vary significantly with various combinations of strength, frequency, and established relations.

### *Pooled relations*

The reliability correlates are a basis for weighting alternative reports. Informants able to provide a more reliable report can be given more weight. We pooled across interviews with the following equation:

$$z_{ij} = \frac{\sum_k (w_{ijk} z_{ijk})}{\sum_k (w_{ijk})},$$

where  $z_{ij}$  is the strength of connection between persons  $i$  and  $j$  in the production network,  $z_{ijk}$  is informant  $k$ 's report on the relation (a score between 0 and 1 in Fig. 4), and  $w_{ijk}$  is a weight indicating whether the informant is in a position to provide a reliable report on the relation:

- = 4.0, if the informant is  $i$  or  $j$  (i.e. the informant is describing one of his or her own relations), or the informant is someone other than  $i$  or  $j$  who has known  $i$  and  $j$  for more than a year, speaks with  $i$  and  $j$  daily, and has more than an average strength relation with  $i$  and  $j$  (first row of Table 4),
- = 0.25, if the informant (a) speaks only weekly with  $i$  and  $j$ , or speaks less than weekly with  $i$  or  $j$ , and (b) has an average or distant connection with  $i$  or  $j$  (bottom two rows of Table 4),
- = 1.0, otherwise (rows two through six of Table 4).

These arbitrary weights capture the three reliability categories in Table 4. The weights highlight the most reliable reports and de-emphasize the least reliable. The most reliable sources of reports have four times, and the least reliable have a fourth, the weight of an average report.

If a relation was captured only once, the weights cancel out in the above equation, and the relation appears in the production network at whatever strength it was reported ( $z_{ij} = z_{ijk}$ ).

But here is an example from the data of a relation captured once, then recaptured four times. Two informants report a strong connection between John and Bob ( $z_{ijk} = 1.0$ ), a third informant reports the connection as average strength ( $z_{ijk} = 0.56$ ), a fourth reports it as distant ( $z_{ijk} = 0.0$ ), and in an interview with John, John himself cites Bob as his most distant connection ( $z_{ijk} = 0.17$ ). The first two informants have average strength relations with John. One is distant from Bob. They speak with John and Bob weekly or less. So, neither informant is in a position to provide a reliable report on the John–Bob relationship ( $w_{ijk} = 0.25$ ). The third informant has strong connections with John and Bob, and has known them for five years, but speaks to them less than weekly. Such a position provides average reliability reports ( $w_{ijk} = 1.0$ ). The fourth informant has strong, frequent, and established relations with John and Bob. John and the fourth informant are in good positions to provide a reliable report on the John–Bob relationship ( $w_{ijk} = 4.0$ ). The pooled relation between John and Bob is defined by the above equation as a weighted average of the five reports;

$$z_{ij} = 0.183$$

$$= (0.25 \times 1.0 + 0.25 \times 1.0 + 1 \times 0.56 + 4 \times 0.0 + 4 \times 0.17) / 9.5,$$

which has a 0.065 variance. Without reliability weights, each report would have equal weight in defining the pooled relation;  $0.546 = (1.0 + 1.0 + 0.56 + 0.0 + 0.17) / 5$ , which has a 0.171 variance.

Two points are illustrated. First, the pooled relation with weighting for reliability is in this example more accurate. It is 0.183, which lies in the distant connection interval of Fig. 4, which is how John and the most reliable informant report the relationship. The unweighted average of 0.546 implies a much stronger relationship. In other words, John and Bob feel distant from one another, but it would be poor form to make a public display of their feelings. They and their closest associates know their distant feelings for one another. To informants farther removed, John and Bob seem close.

Second, there is less error variance in the weighted average. The 0.065 variance across reports weighted for their reliability is less than half the 0.171 unweighted variance. To the extent that less reliable sources of reports create error variance around the pooled estimates

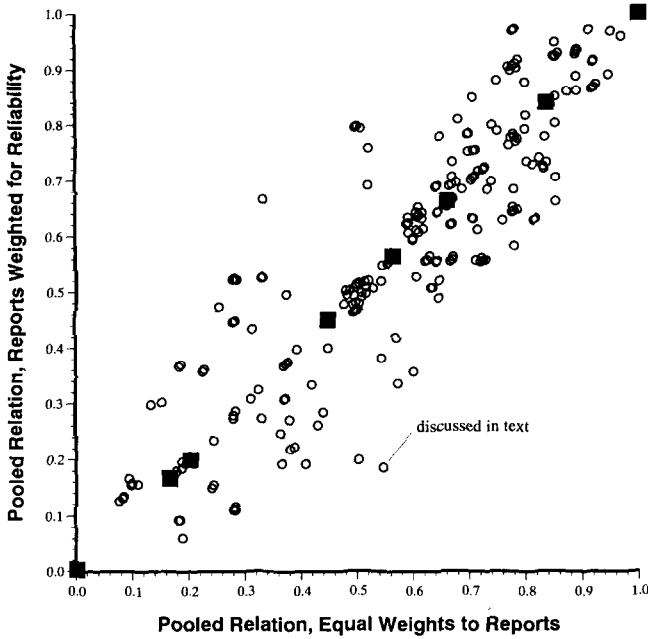


Fig. 8. Pooling reports.

of relations, de-emphasizing less reliable sources lowers the error variance.

Few relations are so affected by the reliability weighting. Figure 8 is a graph of the pooled relations before weighting (horizontal axis) and after (vertical axis). The above relation between John and Bob is marked 'discussed in text' in the graph; 0.55 before weighting, 0.18 after. Few relations lie so far from the diagonal of the graph. The wider spread of relations away from the diagonal for values under 0.6 on the horizontal axis shows that reliability weighting most affects relations that would have been average or distant in the production network.

So how much are the pooled relations improved by the reliability weighting? Very little, or quite a bit; depending on how the question is answered. A quick look at Fig. 8 is enough to see the strong correlation between relations before and after weighting. The exact correlation is 0.99, but most of the relations in the correlation would be unchanged by any method of resolving differences between reports on

the same relation. We captured 2121 relations in the data collection. These are tabulated in Table 1 and plotted in Fig. 8. We have only one report on 1621 of the relations, and informants gave us identical reports on another 175 of the relations (where identical here means identical  $z_{ijk}$  scores, not the broader meaning of Fig. 3). The pooled values of these 1796 relations will be constant across any weighting because they contain no conflicts to resolve. These 1796 relations lie under the solid squares in Fig. 8.

Weighting has its effect on the other 325 relations, relations where we have two or more conflicting reports. These 325 relations are the hollow dots in Fig. 8. Many lie on the diagonal, marking a pooled relation that has the same value before and after weighting. Most lie off the diagonal, like the above discussed relation between John and Bob. The effect of weighting shows up in an analysis of variance in the 325 relations:

Variance between $z_{ij}$	Variance around $z_{ij}$	Percent between
0.0505	0.0571	47%
0.0542	0.0376	59%

The two rows describe relations before and after reliability weighting. The first column is the variance between pooled relations. The hollow dots in Fig. 8 have a 0.0505 variance on the horizontal axis, and a slightly higher 0.0542 variance on the vertical axis. Larger between-relation variance can be expected after weighting to the extent that uninformed reports bias strong and distant relations toward the residual category (like a 'Don't Know' response in this data collection) of average relations. Variance around the pooled relations is error variance. To the extent that informants provide widely different reports on a relationship, we are less certain of the relation's actual strength in the production network. When each report is given equal weight, there is more variance in reports around pooled relations than between them; 47% is between relations. The reliability weighting emphasizes more reliable reports and de-emphasizes reports likely to be unreliable. There is less error variance. The 0.0376 variance around relations with reliability weighting is two-thirds of the 0.0571 variance without weighting.



## 5. Extrapolating to uncaptured relations

We began with 2883 descriptions of 2121 relations. Pooling across recaptures, we have measured the strength of connection in each of the 2121 relationships. That leaves another 22410 relations, in this network of 222 people, on which we have no data. Network analysis involves comparisons across all relations in a network, so we have to assign, implicitly or explicitly, quantitative values to the uncaptured relations. The simplest option is to ignore them. Uncaptured relations can be set to a value of zero. The data collection was designed to capture relations in the production process, so relations not captured are probably disconnections. This assumption's validity depends on the extent to which the fieldwork was sufficient to capture all strong connections. If there is no time to check the assumption, the default is to set uncaptured relations to zero and begin the analysis. Ideally, the assumption can be checked, and found acceptable. We had time to check the assumption. As discussed with Table 1, our fieldwork felt incomplete. Another day or two would have been reassuring. Here are three illustrative indications that our uncaptured relations should not all be set equal to zero.

### *Structure obscured*

Table 5 contains alternative density tables. Five blocks of employees are distinguished; the 32 corporate headquarters employees are distinguished from the four broad functions inside the plant discussed in Fig. 6. Each density is the average relation between people in a row and column. Diagonal densities are mean relations between people in the same function. For example, the 222 person production network includes 19 people who work in the engineering division of the plant. This is the fifth category in Table 5. There are 171 relations between the 19 people, of which 36 relations were captured. The first panel of Table 5 shows that the average relation between engineers is 0.14 if uncaptured relations are set equal to zero (cell 5, 5), the second panel shows that the average value of the 36 captured relationships is 0.64, and the third panel shows that the proportion of relations captured is 0.21 (36 captured over 171 possible).

The first panel is the most generic density table. Densities are low because the captured relations are a small proportion of all relations

Table 5  
Density tables

	Head- quarters	Production Shops	Control	Other	Engi- neering
<i>Mean relations</i>					
Headquarters	0.08				
Shops	0.02	0.08			
Control	0.02	0.05	0.07		
Other	0.01	0.03	0.02	0.04	
Engineering	0.07	0.04	0.05	0.02	0.14
<i>Mean captured relations</i>					
Headquarters	0.58				
Shops	0.34	0.55			
Control	0.38	0.44	0.49		
Other	0.23	0.41	0.40	0.73	
Engineering	0.53	0.43	0.49	0.42	0.64
<i>Proportion relations captured</i>					
Headquarters	0.13				
Shops	0.05	0.14			
Control	0.05	0.11	0.13		
Other	0.03	0.06	0.04	0.05	
Engineering	0.13	0.09	0.10	0.04	0.21

possible. The second panel treats uncaptured relations as missing data. This shows variation in the strength of captured relations, but the densities exaggerate connectivity. Strong connections are more likely to be captured than distant connections, so the mean strength of captured relations is higher than the actual strength of relations. The second panel gives the appearance of strong connections within and between all functions. The third panel depends on the greater tendency for strong connections to be captured. The strength of captured relations is ignored. The question is how many relations are captured. Where many of the possible relations are captured, there must be strong connections. The problem is that many of the captured relations are distant connections. Of the 2121 captured relations, 673 are distant connections in the sense of having a quantitative value under 0.25 (of which 409 have values of 0.00). These distant connections are treated, in the third panel of Table 5, as if they were strong connections.

Two points are illustrated in Table 5. First, setting uncaptured relations to zero obscures much of network structure because so many

relations are uncaptured. The first and third density tables in Table 5 are virtually identical. Densities are about twice as high in the third table, but the stronger connections in the first density table are the stronger connections in the third density table (0.97 correlation between the two tables). The difference between the first and third density tables is only that the strength of captured relations is preserved in the first and ignored in the third. There are so many uncaptured relations that their numbers determine the relative density of connections in the network. The fact that a third of the captured relations (693 of 2121) are distant connections has no effect on network structure at this aggregate level.<sup>6</sup> This is our first indication that uncaptured relations should not all be set equal to zero.

### *Structure ambiguous*

The second point illustrated in Table 5 is substantive. Some qualities are consistent across the three density tables: Density is higher within than between functions. People in production support are the most distant from the rest of the network. The people in engineering have stronger connections to corporate headquarters than people in production, cell (5, 1) is higher than cells (2, 1), (3, 1), and (4, 1).

There are also some important differences. Fig. 9 contains spatial maps. These are multidimensional scalings of the density tables, using Kruskal's (1964) algorithm to preserve metric differences between densities. Functions are close in a map to the extent that the people in them are connected by strong relations. The map at the top of the figure is based on the mean strength of captured and uncaptured relations with uncaptured relations set to zero. This is the first density table in Table 5 (the map for the third density table is, of course,

<sup>6</sup> A practical implication is that if research is intended only to describe aggregate network structure, there is no need to scale or pool the captured relations, just assume that anyone cited is a maximum strength connection, and any two people cited by the same person are connected by a maximum strength relationship. This is the assumption of the third density table in Table 5, and it yields the same structural image as the first density table which is based on variable strength relations with and among cited contacts. If research is intended to describe network structure at a more detailed level, however, variable strengths of relations need to be recorded since so many of the captured relations are in fact distant connections.

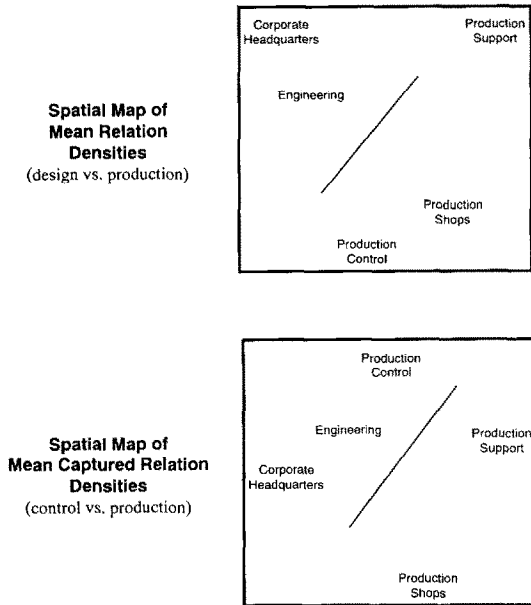


Fig. 9. Spatial maps.

identical). The map at the bottom of Fig. 9 is based on the mean strength of captured relations, the second density table. In both maps, corporate headquarters is at one end and production support is at the other. The maps differ in their primary axis of organization.

According to the mean relation densities, there is a cleavage between corporate headquarters and engineering at the upper left of the map vs. production on the other side. We have drawn a line in the map to indicate the cleavage. This image of the network makes sense because engineering in the plant is incomplete. Much of the design work is done at headquarters, headquarters authorizes certain drawings completed at the plant before the drawings can be used in production, and the head of engineering at the plant reports less to the plant manager than to a senior engineer at corporate headquarters. The basic cleavage in the first map is between design and production.

In the second map, the cleavage is between control and production. Engineering is again close to corporate headquarters, on the opposite side of the map from production support and the production shops.

But the production control function is now connected to engineering and less connected to the actual work of production. Production control appears on the engineering-headquarters side of the cleavage in the network. This image of the network can also make sense. The quality control people work closely with engineering. When defects are identified by quality control, engineering has to certify the defect and define what is needed to correct it. When defects require re-work, delays on other work occur, and the flow of materials is rescheduled, which brings in material control people. The costs involved draw the constant attention of corporate headquarters.

Which is it? Is the plant organized on a contrast between design and production or a contrast between control and production? The point here is not to resolve the question but to note that such a question exists. This is our second indication that the uncaptured relations should not all be set equal to zero. The only additional information we can get at this point is to use the known strengths of the captured relations to make informed guesses about uncaptured relations linking production control to other segments of the production network.

### *Structure incomplete*

Our third indication that the uncaptured should not be set to zero is the rate of change in relations across recapture frequencies. The capture–recapture data collection is a kind of diffusion phenomenon. Strong relations are likely to be recaptured early in the data collection, followed by less strong relations being recaptured as the number of captured relations increases quickly, followed by a period of few new relations being captured as the data collection seeks the remaining, most distant, connections. In other words, the data collection can be expected to generate S-shaped curves like a diffusion process.

Figure 10 displays curves for our data collection. The solid dots to the right of the graph describes what we know about the captured relations. The solid line describes the mean  $z_{ij}$  for relations at each level of recapture. The dashed line describes the proportion of the relations that are distant ( $z_{ij} < 0.25$ ; cf. Fig. 4). The 1621 relations captured only once, for example, have a mean strength of 0.44 and 38% of them are distant. Relations captured multiple times tend to be

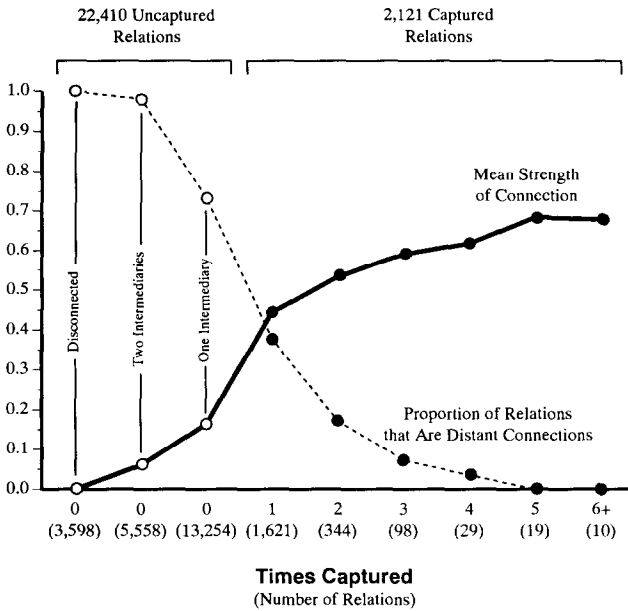


Fig. 10. Extrapolating to uncaptured relations.

stronger. The solid line increases for relations captured twice, three times, and so on. Relations captured multiple times tend not to be distant. The dashed line decreases as recaptures increase. None of the relations captured five or more times are distant connections. Few relations were captured more than six times, so the final column in the graph is six or more captures (cf. Table 1).

The hollow dots to the left of the graph describes our best guesses about uncaptured relations. If we did not capture the relation between Bill and John, we looked through the relations we captured for some third party whose relations with Bill and John are known. The third party is an intermediary that completes a connection between Bill and John. The indirect connection could be strong, if the third party is close to Bill and John, or it could be distant, if the third party is distant from either Bill or John. Bill's relation with the third party times the third party's relation with John is a measure of connection between Bill and John. For each uncaptured relation  $z_{ij}$ , we found every third party  $k$  whose relation with  $i$  and with  $j$  was captured,  $z_{ij}$

is the average value of  $z_{ik}z_{kj}$ .<sup>7</sup> These are the hollow dots in Fig. 10. Of the uncaptured relations, more than half can be filled with an indirect connection through one intermediary (13 254 of the 22 410). Of the more distant remaining relations, more than half can be filled with an indirect connection through two intermediaries (5558 of the 9156). The remaining 3598 uncaptured relations can not be filled through any number of intermediaries.

Three points are illustrated. First, we captured the core, or spine, of the production network in the sense that the many uncaptured relations are quickly filled in with the relations we did capture. Almost all of the uncaptured relations can be filled with indirect connections through one or two intermediaries (18 812 is 84% of the 22 410 uncaptured relations). The rest are likely to be disconnections since there are no combinations of connections through the core of the network that can fill them.

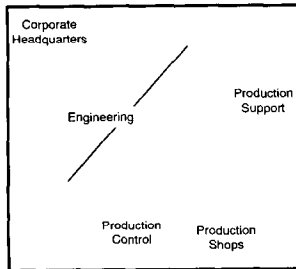
Second, there are many average and strong connections among the uncaptured relations. The hollow dots in Fig. 10 show the increasing tendency for indirect connections to be distant, but the large number of relations involved means that many new connections are being made. Of the indirect relations through one intermediary, a large proportion, 73%, are distant connections. But the remaining 27% amount to 3556 relations, 3394 average strength connections with a value from 0.25 to 0.75, and another 162 strong connections with a value over 0.75. This is a substantial addition to the 1002 average strength relations and 446 strong relations that were captured.

<sup>7</sup> The path distance analogue would be to find the person  $k$  most strongly connected to  $i$  and  $j$ , then set the uncaptured  $z_{ij}$  equal to the maximum  $z_{ik}z_{kj}$ . This was our first choice for assigning values to the uncaptured relations. To check the validity of the imputation, we computed maximum  $z_{ik}z_{kj}$  for the relations we had captured as well as those we had not captured. If the imputation is valid, imputed scores should be strongly correlated with the known scores. The correlation is only 0.32 between captured  $z_{ij}$  and maximum  $z_{ik}z_{kj}$ . The problem is that most pairs of people in the network are strongly connected to some third party, so maximum  $z_{ik}z_{kj}$  overstates the strength of connection. It obscures disconnections in the network. Among the 673 distant connections, the average indirect connection is strong (mean maximum  $z_{ik}z_{kj} = 0.61$ ), even though we know from the informants that these are distant connections (mean  $z_{ij} = 0.06$ ). Instead of equating uncaptured relations to the strongest indirect connection, we set them equal to the typical indirect connections, the mean of captured relations  $z_{ik}z_{kj}$ . The resulting indirect connections are more consistent with the known strengths of relations, so we have more confidence in them as a measure of uncaptured relations. In contrast to the 0.32 correlation between captured  $z_{ij}$  and the maximum  $z_{ik}z_{kj}$ , there is a 0.70 correlation between  $z_{ij}$  and the mean of captured  $z_{ik}z_{kj}$ .

Third, the S-shape diffusion curves are evident in Fig. 10, and provide a quick indication of data collection's progress. Fieldwork is nearly complete when the derivative of change in relations across recapture frequency reverses sign. The dashed line in Fig. 10 describes distant connections at each recapture frequency:

recapture frequency	distant connections	cumulative
6 +	0	0
5	0	0
4	1	1
3	7	8
2	59	67
1	606	673
one intermediary	9698	10371
two intermediaries	5440	15811
disconnected	3598	19409

**Spatial Map of Mean Relation Densities**  
(design vs. production)



**Spatial Map of Mean Captured Relation Densities**  
(design vs. production)

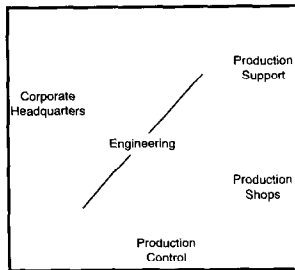


Fig. 11. Spatial maps considering indirect connections.



Among the most recaptured relations, none are distant. One of the relations captured four times is distant, seven of the three-capture relations, and so on. The rate of capturing new distant connections is accelerating. The rate decelerates at one intermediary connection (derivative reverses sign), is slower at two intermediaries, and is flat thereafter. The rate of acquiring distant connections is still accelerating among our once-captured relations, so the fieldwork is incomplete. It would have been complete if the rate at which we were acquiring new distant connections had begun to slow. In this study population, judging from Fig. 10, that would have been when 70% or so of our once-captured relations were distant. The percentage of distant connections is not the key variable. That will vary with the social structure of different study populations. The key variable is the rate at which new distant connections are being acquired when the fieldwork ends.<sup>8</sup>

<sup>8</sup> Diffusion models can be used to estimate the number of disconnections. Let  $\mathbb{N}$  be the number of distant connections in the network. Where recapture frequencies are listed in descending order as on page 130, let  $N(t)$  be the cumulative number of distant connections captured at recapture level  $t$ ;  $N(0)$  equals 0 for relations captured five times,  $N(1)$  equals 1 for relations captured four times,  $N(2)$  equals 8 for relations captured three times, and so on. Change in the cumulative number of captured distant connections can be described by a diffusion model (Bass 1969; Mahajan and Wind 1986):  $dN(t)/dt = \beta N(t)[\mathbb{N} - N(t)]$ , where  $\beta$  is a contagion effect coefficient,  $N(t)$  is the number captured, and  $[\mathbb{N} - N(t)]$  is the number that remain uncaptured. This is the metric version of the contagion model familiar to network analysts from Coleman *et al.*'s (1966) *Medical Innovation Study* (see Burt 1987: 1271n):  $dy/dt = \beta y(1 - y)$ , where  $y$  is the cumulative proportion of adoptions at time  $t$  ( $y = N(t)/\mathbb{N}$ ), and the following difference equation:

$$N(t+1) - N(t) = \beta \mathbb{N} N(t) - \beta N(t)^2 = \alpha N(t) - \beta N(t)^2,$$

can be used to estimate  $\beta$  directly, and  $\mathbb{N}$  thereafter ( $\mathbb{N} = \alpha/\beta$ ). From the rate at which distant connections are captured, in other words, the model looks for a diffusion S-shaped curve that fits the data and predicts the number of distant connections in the network from the ceiling in the S-curve. Ordinary least-squares regression across seven recapture levels yields values of 1.61100 for  $\alpha$  and 0.000090 for  $\beta$ , which yields an estimate of 17900 distant connections in the network. Direct estimation with a nonlinear least-squares algorithm yields an estimate of 0.000090 for  $\beta$ , and 17872 distant connections in the network. A more sophisticated model with a third parameter describing the constant tendency for a distant connection to be captured (corresponding to personal preference in diffusion) yields a slightly higher 18041 distant connections (and assured us that the third parameter is negligible; 0.2  $t$ -test). In short, our 19409 distant connections is probably too high. The true number is closer to 18000. However, we do not depend on a binary distinction between average and distant (many of our distant connections have quantitative strengths close to the 0.25 cut-off for average strength connections), and the predicted 73% distant connections in the network (18000 of 24531 connections) is very close to our observed 79% (19409 of 24531 connections).

Our conclusion is that we should use indirect connections as our best guess about the strength of uncaptured relations. It is clear from Fig. 10 that the fieldwork was incomplete, but the captured relations are well-located in the network to fill in relations that were not captured. Also, the problem of distant connections having no effect on network structure at the aggregate level is resolved. Recall the 0.97 correlation between the first and third density tables in Table 5. The correlation is now only 0.23 between mean relations and proportion captured relations (where captured now includes indirect connections), because we have captured a high proportion of relations in each block of the density table. Also, the structural ambiguity illustrated in Fig. 9 is resolved. The spatial maps in Fig. 11 are the same as the maps in Fig. 9 except densities for the maps in Fig. 11 include indirect connections. The two maps present a similar image of the production network. Engineering is the link between corporate headquarters on one side and production on the other.

## **6. Summary**

Based on the methodological results presented, and substantive analysis (not reported) of the final 222 person production network, we conclude that the capture–recapture strategy can be a useful way to measure a large network quickly. We have discussed five stages to the work:

(1) Conduct survey network interviews with people (informants) positioned in the study population such that their contacts overlap to provide recaptured relations. For the example discussed here, this first stage required a day with an informant to become familiar with the business and identify people to interview, then three days of interviews. Forty-five informants were interviewed. Their responses defined 2121 relations in a network of 222 people (Table 1).

(2) Determine reliability from data consistency across recaptures. With segments of the network based on lean data, data reliability is especially important. We had 1210 capture–recapture pairs in which one informant’s report on a relationship could be compared to a second informant’s report on the same relationship. These paired reports are most often identical and rarely contradictory (Fig. 2).

Informants were as reliable in describing relations between other people as they were about their own relations with other people (Fig. 3).

(3) Triangulate network response categories to assign quantitative scores to levels of relationship. To pool recaptured relations across interviews and model the network, we needed to convert the response categories describing qualitative levels of connection strength into quantitative scores. We used a one-dimension loglinear association model fit to capture–recapture pairs and balance pairs (Table 2 and Fig. 4). Response bias can distort the balance pairs. We describe a *tertius* bias in which informants exaggerate the extent to which they broker the connection between their closest contacts (Figs. 5 and 6). This might be a response bias peculiar to business networks, but we think not.

(4) Use reliability correlates to weight recaptured relations to compute the final network pooled across interviews. To reduce error around the pooled relations, we weighted reports for reliability. This required other data on the relations and knowing how reliability was correlated with the other data. Strong connections in the plant are rare between people who meet less than once a week or who have known one another for less than a year. A year acquaintance establishes a relationship, after which connections develop as strong or distant regardless of the years for which two people are acquainted. Strong relations tend to involve daily contact. Reliability varies with an informant's ties to the people connected by the relationship the informant is describing. The aggregate tendency for paired reports to be identical is significantly higher when both informants have strong-frequent-established relations with the contacts between whom they are describing a relationship (Table 4). Reliability weighting emphasizes more reliable reports and gives less weight to reports likely to be unreliable. The 0.0376 error variance with reliability weighting is two-thirds of the 0.0571 error variance without weighting.

(5) Extrapolate from the known strengths of the captured relations to define the uncaptured relations. We captured only 2121 of the 24 531 relations in our network of 222 people. The captured relations were clearly the core, or spine, of the network in the sense that over half of the uncaptured relations could be completed through captured relations with one intermediary, and over half of the remaining, more distant, uncaptured relations could be completed through captured

relations with two intermediaries. The final network has a global density of 0.14 with extremes of 19 409 distant connections and 616 strong connections.

## References

- Bass, Frank M.  
1969 "A new product growth model for consumer durables". *Management Science* 15: 215–227.
- Bernard, H. Russell, Eugene C. Johnson, Peter D. Killworth and Scott Robinson  
1989 "Estimating the size of an average personal network and of an event subpopulation". In: Manfred Kochen (Editor) *The Small World*. Norwood, NJ: Ablex, pp. 159–175.
- Burt, Ronald S.  
1987 "Social contagion and innovation, cohesion versus structural equivalence". *American Journal of Sociology* 92: 1287–1335.  
1992 *Structural Holes: The Social Structure of Competition*. Cambridge: Harvard University Press.
- Burt, Ronald S. and Miguel G. Guilarte  
1986 "A note on scaling the General Social Survey network item response categories". *Social Networks* 8: 387–396.
- Burt, Ronald S. and Michael J. Minor (Editors)  
1983 *Applied Network Analysis*. Beverly Hills: Sage Publications.
- Burt, Ronald S. and Don Ronchi  
1990 "Contested control in a large manufacturing plant". In: J. Weesie and H. Flap (Editors) *Social Networks through Time*. Utrecht, Holland: ISOR, University of Utrecht, pp. 121–157.
- Coleman, James S.  
1958 "Relational analysis: the study of social organizations with survey methods". *Human Organization* 16: 28–36.
- Coleman, James S., Elihu Katz and Herbert Menzel  
1966 *Medical Innovation*. New York: Bobbs-Merrill.
- Erickson, Bonnie H.  
1978 "Some problems of inference from chain data". In: Karl F. Schuessler (Editor) *Sociological Methodology 1979*. San Francisco: Jossey-Bass, pp. 276–302.
- Frank, Ove  
1978 "Sampling and estimation in large social networks". *Social Networks* 1: 91–101.
- Freeman, Linton C., Sue C. Freeman and Alaina G. Michaelson  
1989 "How humans see social groups: a best of the Sailer–Gaulin models". *Journal of Quantitative Anthropology* 1: 229–238.
- Freeman, Linton C. and Claire R. Thompson  
1989 "Estimating acquaintanceship volume". In: Manfred Kochen (Editor) *The Small World*. Norwood, NJ: Ablex, pp. 147–158.
- Faust, Katherine and Stanley Wasserman  
1993 "Correlation and association models for studying measurements on ordinal relations," In: Peter V. Marsden (Editor) *Sociological Methodology 1993*. Cambridge, MA: Basil Blackwell.
- Goodman, Leo A.  
1961 "Snowball sampling". *Annals of Mathematical Statistics* 34: 148–170.

- 1984 *The Analysis of Cross-Classified Data Having Ordered Categories*. Cambridge: Harvard University Press.
- Granovetter, Mark S.  
1976 "Network sampling: some first steps". *American Journal of Sociology* 81: 1287–1303.
- Johnson, Jeffrey C.  
1990 *Selecting Ethnographic Informants*. Newbury Park, CA: Sage Publications.
- Johnson, Jeffrey C., James S. Boster and D. Holbert  
1989 "Estimating relational attributes from snowball samples through simulation". *Social Networks* 11: 135–158.
- Klovdahl, Alden S.  
1989 "Urban social networks: some methodological problems and possibilities". In: Manfred Kochen (Editor) *The Small World*. Norwood, NJ: Ablex, pp. 176–210.
- Kruskal, Joseph B.  
1964 "Multidimensional scaling by optimizing goodness of fit to a nonmetric hypothesis." *Psychometrika* 29: 1–27.
- Mahajan, Vijay and Yoram Wind (Editors)  
1986 *Innovation Diffusion Models of New Product Acceptance*. Cambridge, MA: Ballinger Publishing.
- Marsden, Peter V.  
1990 "Network data and measurement". *Annual Review of Sociology* 16: 435–463.
- Pool, Ithiel de Sola, and Manfred Kochen  
1978 "Contacts and influence". *Social Networks* 1: 5–51.
- Romney, A. Kimball, and Susan C. Weller  
1984 "Predicting informant accuracy from patterns of recall among individuals". *Social Networks* 6: 59–77.
- Snijders, Tom A.B.  
1992 "Estimation on the basis of snowball samples: how to weight?" Paper presented at the "Workshop on Generalizability Questions for Snowball Sampling and Other Ascending Methodologies," University of Groningen, Groningen, Netherlands.