

---

# The Occurrence of Singularities in Cosmology

S. W. Hawking

*Proc. R. Soc. Lond. A* 1966 **294**, 511-521

doi: 10.1098/rspa.1966.0221

---

## Email alerting service

Receive free email alerts when new articles cite this article - sign up in the box at the top right-hand corner of the article or click [here](#)

---

To subscribe to *Proc. R. Soc. Lond. A* go to: <http://rspa.royalsocietypublishing.org/subscriptions>

---

# The occurrence of singularities in cosmology

BY S. W. HAWKING

*Gonville and Caius College, University of Cambridge*

*(Communicated by H. Bondi, F.R.S.—Received 15 April 1966)*

It is shown that singularities of space-time are inevitable if the Einstein equations hold, if matter has normal properties and if the universe satisfies certain reasonable global conditions. The singularities would be in the past and would, in principle, be observable. Observation to determine whether such singularities actually occurred would provide a powerful test of the Einstein equations in strong fields. The singularity would not necessarily constitute a beginning of the universe.

## 1. INTRODUCTION

If it is assumed that there is no continual creation of matter, the observed recession of distant galaxies indicates that the density of the universe must have been higher in the past. There are two main possibilities: the universe might have contracted to some maximum density and then expanded again, or there might have been a time when the density was infinite. This would imply a singularity of space-time. Robertson (1933) showed that the universe must have had a singularity if it was exactly spatially homogeneous and isotropic, if the Einstein equations without cosmological constant held, and if matter had normal properties. However, it was suggested that maybe the singularity only occurred because of the high symmetry implied by homogeneity and isotropy and that a more realistic model with fewer symmetries would not have a singularity. Raychaudhuri (1955) and Komar (1956) showed that there must be a singularity in a model containing non-rotating dust and Shepley (1964) and Hawking & Ellis (1965) showed the same for a model that was spatially homogeneous but anisotropic. However, these models are not fully general, since in all of them there is some exact restriction or symmetry on the flow of matter. Lifshitz & Khalatnikov (1963) claim to have proved that there will not be a singularity in a fully general model that is, one that has the full number of arbitrary functions. However Penrose (1965) has shown that if the universe possesses a noncompact Cauchy surface a collapsing star must develop a singularity. The same method may be applied to prove the occurrence of a singularity in a universe model with a noncompact Cauchy surface which is approximately homogeneous and isotropic (Hawking 1965). However, the Cauchy surface of the universe might be compact (a 'closed' model) or it might even be that the universe possessed no global Cauchy surface. This paper will consider these possibilities and it will be shown that under certain reasonable conditions it can still be proved that singularities must occur.

## 2. NOTATION

Space-time is represented as a connected paracompact four-dimensional  $C^\infty$  manifold  $\mathcal{M}$  with pseudo-Riemannian metric tensor  $g_{ij}$  of signature minus two. Covariant differentiation is indicated by a semi-colon, covariant differentiation along a timelike line by  $D/Ds$  and Lie differentiation by  $\mathcal{L}$ . Square brackets around indices indicate antisymmetrization. The conventions for the Riemann and Ricci tensors are

$$\begin{aligned} 2V_{a;[bc]} &= R^p{}_{abc}v_p, \\ R_{ab} &= R^c{}_{acb}. \end{aligned}$$

We assume the Einstein equations,

$$R_{ab} - \frac{1}{2}g_{ab}R = T_{ab},$$

where  $T_{ab}$  is the energy-momentum tensor of matter and units are such that the speed of light and the constant of gravitation equal one.

Space-time will be said to be singularity-free if the metric is a field of class  $C^2$  and  $\mathcal{M}$  is geodesically complete with respect to the metric.

## 3. CONJUGATE POINTS

A timelike congruence will be said to be geodesic irrotational if its unit tangent vector  $v^a$  satisfies

$$v_{[a;b]} = 0.$$

The expansion of such a congruence is defined as  $\theta = v^a{}_{;a}$ . A point  $p$  will be said to be a singular point on a geodesic  $\gamma$  of a congruence if  $\theta$  is infinite on  $\gamma$  at  $p$ . A point  $p$  will be said to be conjugate to a point  $q$  along a geodesic  $\gamma$  from  $q$  to  $p$  if  $p$  is a singular point on  $\gamma$  of the irrotational congruence of all timelike geodesics through  $q$ . A point  $p$  will be said to be conjugate to a spacelike hypersurface  $\mathcal{H}$  if  $p$  is a singular point of the irrotational congruence of geodesics normal to  $\mathcal{H}$ .

A more precise description of conjugate points is as follows: let  $k^a$  be a vector connecting points corresponding distances along neighbouring geodesics in a congruence. Then  $k^a$  is dragged along by the congruence, that is,

$$\mathcal{L}_{v^a} k^a = 0. \quad (1)$$

Therefore

$$Dk^a/Ds = v_{a;b}k^b \quad (2)$$

and

$$D^2k^a/Ds^2 = R_a{}^b{}_{bc}v^a v^c k^b. \quad (3)$$

Introducing an orthonormal tetrad  $e^a(m = 1, 2, 3, 4)$  which is parallelly transported along the geodesic with  $e^a = v^a$ , equation (3) may be written

$$\frac{d^2k}{ds^2} = a \frac{k}{nr}, \quad (4)$$

where

$$\begin{aligned} k &= e^a k_a, \\ a &= e^b e^c R_{abcd} v^a v^d. \end{aligned}$$

A solution of equation (4) will be called a Jacobi field. There are eight independent solutions. Since  $v^a$  and  $sv^a$  are solutions, the other six solutions may be taken orthogonal to  $v^a$ .

If  $q$  is a point on a geodesic  $\gamma$  with tangent vector  $v^a$ , the Jacobi fields along  $\gamma$  orthogonal to  $v^a$ , which vanish at  $q$  may be regarded as generating neighbouring geodesics in the irrotational congruence of all geodesics through  $q$ . Thus they obey

$$dk_m/ds = v \underset{m:n}{k}^n, \quad (5)$$

where

$$v \underset{m:n}{k}^n = e^a e^b v_{a;b} \underset{m:n}{k}^n \quad (m, n = 1, 2, 3).$$

They may be written

$$k \underset{m}{n} = A \underset{mn}{(s)} \frac{dk}{ds} \Big|_q,$$

where  $A \underset{mn}{(s)}$  is a  $3 \times 3$  matrix which relates the fields at any point on  $\gamma$  to their derivatives at  $q$ .

For the rest of this paper we will omit tetrad indices from one to three and write product in matrix form, e.g.

$$A \underset{mn}{k}^n = \mathbf{A} \cdot \mathbf{k}.$$

Note that because of the signature of the metric

$$\mathbf{A} \cdot \mathbf{k} = - \sum_n A \underset{mn}{k}^n.$$

Then by (5)

$$d\mathbf{A}/ds = \mathbf{V} \cdot \mathbf{A},$$

where

$$(\mathbf{V})_{mn} = v \underset{m:n}{k}^n.$$

Thus

$$\begin{aligned} \theta = \mathcal{F}^*(\mathbf{V}) &= \mathcal{F}^* \left( \frac{d\mathbf{A}}{ds} \cdot \mathbf{A}^{-1} \right) \\ &= [\det(\mathbf{A})]^{-1} \frac{d}{ds} [\det(\mathbf{A})]. \end{aligned}$$

But

$$d^2\mathbf{A}/ds^2 = \mathbf{a} \cdot \mathbf{A}.$$

Therefore  $d\mathbf{A}/ds$  is finite. Hence  $\theta$  is infinite where and only where  $\det(\mathbf{A})$  is zero. But  $\det(\mathbf{A})$  is zero if and only if there is a nonzero  $\mathbf{h}$  such that

$$\mathbf{A} \cdot \mathbf{h} = 0.$$

This gives an alternative definition of conjugate points.  $p$  is conjugate to  $q$  along  $\gamma$  if and only if there is a Jacobi field along  $\gamma$ , not identically zero which vanishes at  $p$  and  $q$ . This also shows that singular points of congruences are points where infinitesimally neighbouring geodesics meet.

## 4. THE VARIATION FORMULAE

The length  $L$  of a timelike path  $\gamma$  from  $q$  to  $p$  is defined as

$$L = \int_q^p ds.$$

We wish to investigate the behaviour of  $L$  under a variation of the path  $\gamma$  from  $q$  to  $p$ . A variation  $\alpha$  of  $\gamma$  keeping end-points fixed is defined as a continuous map

$$\alpha: (-\epsilon, \epsilon) \times [0, L] \rightarrow \mathcal{M}$$

for some  $\epsilon > 0$ , such that

- (1)  $\alpha(0, s) = \gamma(s)$ ;
- (2) there is a subdivision  $0 = s_1 < s_2 \dots < s_k = L$  of  $[0, L]$  such that  $\alpha$  is  $C^\infty$  on each strip  $(-\epsilon, \epsilon) \times [s_i, s_{i+1}]$ ;
- (3) the variation vector  $w^a = \partial\alpha(u, s)/\partial u|_{u=0}$  is orthogonal to the tangent vector  $v^a = \partial\alpha(0, s)/\partial s$ ;
- (4)  $\alpha(u, 0) = q$ ,  $\alpha(u, L) = p$ .

The derivative of  $L$  with respect to the variation  $\alpha$  is

$$\left. \frac{\partial L(\alpha(u))}{\partial u} \right|_{u=0} = -\Sigma w_a [v^a] - \int_0^L w_a \frac{Dv^a}{Ds} ds,$$

where  $[v^a]$  is discontinuity in  $v^a$ . Thus  $\partial L/\partial u$  is zero under all variations if and only if  $Dv^a/Ds = 0$  and there are no discontinuities in  $v^a$ . This expresses the well known fact that a geodesic (i.e. a path with  $Dv^a/Ds = 0$ ) has stationary length. For a geodesic we may investigate the second derivative of  $L$  under a two parameter variation,

$$\alpha: U \times [0, L] \rightarrow \mathcal{M},$$

where  $U$  is a neighbourhood of  $(0, 0)$  in  $E^2$ , with properties as before and

$$w_1^a = \left. \frac{\partial\alpha(u_1, 0, s)}{\partial u_1} \right|_{u_1=0},$$

$$w_2^a = \left. \frac{\partial\alpha(0, u_2, s)}{\partial u_2} \right|_{u_2=0}$$

are the two variation vectors.

Then the second derivative of  $L$  is

$$\left. \frac{\partial^2 L(\alpha(u_1, u_2))}{\partial u_1 \partial u_2} \right|_{\substack{u_1=0 \\ u_2=0}} = - \int_0^L w_{2a} \left( \frac{D^2 w_1^a}{Ds^2} + R^a{}_{bcd} v^b v^c w_1^d \right) ds - \Sigma w_{2a} \left[ \frac{Dw_1^a}{Ds} \right]$$

(for derivation see Milnor 1963, p. 75).

We may regard this second derivative of  $L$  as a bilinear functional  $I(\mathbf{w}_1, \mathbf{w}_2)$  of the space of piecewise  $C^\infty$  vector fields along  $\gamma$ , orthogonal to  $v^a$  and vanishing at  $q$  and  $p$ .

LEMMA 1. *If  $q$  has no conjugate points on  $\gamma$  in  $qp$  then  $I(\mathbf{w}_1, \mathbf{w}_2)$  is negative definite, and therefore  $\gamma$  is maximal under variations  $\alpha$ .*

Consider a set of three variation fields along  $\gamma$  orthogonal to  $v^a$  and vanishing at  $q$  and  $p$ . They may be written in the form

$$\mathbf{w}_i = \mathbf{C}(s) \cdot \mathbf{l}_i \quad (i = 1, 2, 3),$$

where  $\mathbf{l}_i$  are a basis orthogonal to  $v^a$  at  $q$ . Let us take  $d\mathbf{w}_i/ds$  to be linearly independent at  $p$ . The matrix  $\mathbf{C}(s)$  vanishes at  $q$  and  $p$ . As there are no points conjugate to  $q$  in  $qp$ , the matrix  $\mathbf{A}(s)$  of the previous section vanishes at  $q$  and is nonsingular elsewhere in  $qp$ . Thus  $\mathbf{C}(s)$  may be written as

$$\mathbf{C}(s) = \mathbf{A}(s) \cdot \mathbf{B}(s),$$

where  $\mathbf{B}(s)$  vanishes at  $p$ .

Then

$$\begin{aligned} I(\mathbf{w}_i, \mathbf{w}_j) &= \mathbf{l}_i^T \cdot \left\{ -\Sigma \mathbf{B}^T \cdot \mathbf{A}^T \cdot \left[ \frac{d}{ds} (\mathbf{A} \cdot \mathbf{B}) \right] - \int_0^L \mathbf{B}^T \cdot \mathbf{A}^T \cdot \left( \frac{d^2}{ds^2} (\mathbf{A} \cdot \mathbf{B}) - \mathbf{a} \cdot \mathbf{A} \cdot \mathbf{B} \right) ds \right\} \cdot \mathbf{l}_j \\ &= \mathbf{l}_i^T \cdot \left\{ -\Sigma \mathbf{B}^T \cdot \mathbf{A}^T \cdot \left[ \frac{d}{ds} (\mathbf{A} \cdot \mathbf{B}) \right] - \int_0^L \mathbf{B}^T \cdot \mathbf{A}^T \cdot \left( \mathbf{A} \cdot \frac{d^2 \mathbf{B}}{ds^2} + 2 \frac{d\mathbf{A}}{ds} \cdot \frac{d\mathbf{B}}{ds} \right) ds \right\} \cdot \mathbf{l}_j \\ &= \mathbf{l}_i^T \cdot \left\{ \int_0^L \left( \frac{d\mathbf{B}^T}{ds} \cdot \mathbf{A}^T \cdot \mathbf{A} \cdot \frac{d\mathbf{B}}{ds} + \mathbf{B}^T \cdot \left( \frac{d\mathbf{A}^T}{ds} \cdot \mathbf{A} - \mathbf{A}^T \cdot \frac{d\mathbf{A}}{ds} \right) \cdot \frac{d\mathbf{B}}{ds} \right) ds \right\} \cdot \mathbf{l}_j. \end{aligned}$$

But  $d\mathbf{A}/ds = \mathbf{V} \cdot \mathbf{A}$ , where  $\mathbf{V}$  is symmetric. Therefore,

$$I(\mathbf{w}_i, \mathbf{w}_j) = \mathbf{l}_i^T \cdot \left\{ \int_0^L \frac{d\mathbf{B}^T}{ds} \cdot \mathbf{A}^T \cdot \mathbf{A} \cdot \frac{d\mathbf{B}}{ds} ds \right\} \cdot \mathbf{l}_j.$$

But  $(d\mathbf{B}^T/ds) \mathbf{A}^T \mathbf{A} (d\mathbf{B}/ds)$  is negative semi-definite (because of the signature of the metric), and  $d\mathbf{B}/ds$  is nonsingular at  $p$ . Therefore  $I(\mathbf{w}_i, \mathbf{w}_j)$  is negative definite.

LEMMA 2. *If there is a point  $r$  on  $\gamma$  in  $qp$ , conjugate to  $q$  then  $\gamma$  is not maximal from  $p$  to  $q$ .*

Let  $\mathbf{k}$  be a Jacobi field along  $\gamma$  vanishing at  $q$  and  $r$ . Let  $\mathbf{h}$  be any field along  $\gamma$  vanishing at  $q$  and  $p$  and satisfying  $\mathbf{h}^T \cdot d\mathbf{k}/ds = 1$  at  $r$ . Extend  $\mathbf{k}$  to  $p$  by putting it zero in  $rp$ . Let  $\mathbf{w}$  be the field,

$$\mathbf{w} = c\mathbf{h} + c^{-1}\mathbf{k},$$

where  $c$  is a constant. Then

$$I(\mathbf{w}, \mathbf{w}) = 2 - c^2b,$$

where  $b$  is a constant independent of  $c$ . Thus by taking  $c$  small enough  $I(\mathbf{w}, \mathbf{w})$  can be made positive. Thus a geodesic  $\gamma$  from  $q$  to  $p$  is maximal if and only if  $q$  has no conjugate points in  $qp$ .

We may also consider variations of  $\gamma$  in which the end point  $q$  is not kept fixed but is allowed to vary over a surface  $\mathcal{H}$  orthogonal to  $\gamma$ . Consider variations in  $L$  where  $L$  is the distance from  $\mathcal{H}$  to  $p$ . Then

$$\begin{aligned} I(\mathbf{w}_1, \mathbf{w}_2) &= -\Sigma \mathbf{w}_1^T \cdot \left[ \frac{d\mathbf{w}_2}{ds} \right] - \int_0^L \mathbf{w}_1^T \cdot \left( \frac{d^2 \mathbf{w}_2}{ds^2} - \mathbf{a} \cdot \mathbf{w}_2 \right) ds \\ &\quad - \left( \mathbf{w}_2^T \cdot \frac{d\mathbf{w}_1}{ds} - \mathbf{w}_2^T \cdot \mathbf{U} \cdot \mathbf{w}_1 \right) \Big|_q, \end{aligned}$$

where  $(\mathbf{U})_{mn} = u_{m:n}$ , and  $u^a$  is the unit normal to  $\mathcal{H}$ .

LEMMA 3. *A geodesic  $\gamma$ , orthogonal to  $\mathcal{H}$ , from  $q$  on  $\mathcal{H}$  to  $p$  is maximal if and only if there is no point conjugate to  $\mathcal{H}$  in  $qp$ .*

The proof is similar to that for lemmas 1 and 2.

### 5. THE ENERGY CONDITION

Equations of state are idealizations that are only approximately obeyed by the actual matter. Thus a singularity whose occurrence was dependent on the exact form of a particular equation of state would not be physically realistic. We shall therefore not assume anything about the detailed nature of the energy-momentum tensor  $T_{ab}$  but merely certain physically reasonable inequalities. The energy density in the frame of an observer with velocity  $u^a$  is  $E = T_{ab}u^au^b$ . We shall assume that this is greater than or equal to  $\frac{1}{2}T^a_a$  for any velocity,  $u^a$ . For a fluid with an isotropic pressure this assumption holds if the density is nonnegative and the pressure is greater than minus one-third the density. These are properties that any normal matter should have.

The conditions we have assumed imply  $R_{ab}v^av^b \geq 0$  for any timelike unit vector  $v^a$ . The significance of this is as follows: for any geodesic irrotational timelike congruence, the expansion  $\theta$  obeys Raychaudhuri's equation,

$$\theta_{;b}v^b = -v_{a;b}v^{b;a} - R_{ab}v^av^b.$$

Since  $v_{[a;b]} = 0$ ,  $v_{a;b}v^{b;a} \geq \frac{1}{3}\theta^2$ . Therefore if  $\theta$  is negative it will go to minus infinity in a finite distance. That is to say a congruence that is converging will have a singular point within a finite distance where neighbouring geodesics intersect. This result will be used in the next section.

### 6. SINGULARITIES

THEOREM 1. *Space-time cannot be singularity-free if,*

- (1) *The energy-momentum tensor satisfies the condition of the previous section.*
- (2) *The null-cones can be divided continuously into two classes 'past' and 'future'. (This means we can globally assign a time direction.)*
- (3) *There is a Cauchy surface  $\mathcal{H}$ . (We define a Cauchy surface as a complete connected spacelike  $C^\infty$  surface which intersects every timelike and null line once and once only.)*

(4) *The expansion of  $v^a$ , the unit normal to  $\mathcal{H}$ , has a positive lower bound  $c$  on  $\mathcal{H}$ , i.e.  $v^a_{;a} \geq c > 0$ .*

Condition (2) is necessary if our ideas of time are valid. If condition (3) holds, condition (4) might be regarded as a precise statement of what we mean when we say the universe is expanding. Condition (3) will be further discussed in the next section.

To prove the theorem we will assume that space-time is singularity-free and show that this leads to an inconsistency if conditions (1) to (4) hold. The main lines of the proof will be given. It is necessary first to establish several lemmas.

Let  $\mathcal{F}(\mathcal{S})$  and  $\mathcal{P}(\mathcal{S})$  be the points that can be reached from a set  $\mathcal{S}$  by respectively future and past directed timelike line. They have the property that if  $p \in \mathcal{P}(\mathcal{S})$  then  $\mathcal{P}(p) \subset \mathcal{P}(\mathcal{S})$ .  $\mathcal{P}(\mathcal{S})$  and  $\mathcal{F}(\mathcal{S})$  are open sets.

The boundary of the set  $\mathcal{S}$  will be denoted by  $\dot{\mathcal{S}} \equiv \overline{\mathcal{S}} \cap \overline{(\mathcal{M} - \mathcal{S})}$ , where a bar indicates closure.

**LEMMA 4.** *If a nonempty set  $\mathcal{N}$  is such that  $p \in \mathcal{N}$  implies  $\mathcal{P}(p) \subset \mathcal{N}$  and if  $\mathcal{M} - \mathcal{N}$  is not empty then  $\mathcal{N}$  is locally homeomorphic to  $E^3$  and no two points of  $\mathcal{N}$  have timelike separation.*

If  $\mathcal{N}$  and  $\mathcal{M} - \mathcal{N}$  are nonempty, then  $\dot{\mathcal{N}}$  is non-empty. Let  $q \in \dot{\mathcal{N}}$ , then  $\mathcal{P}(q) \subset \mathcal{N}$  and  $\mathcal{F}(q) \subset \mathcal{M} - \mathcal{N}$ . Thus there can be no point  $r \in \dot{\mathcal{N}} \cap \mathcal{F}(q)$  since if  $r \in \mathcal{F}(q)$  then there would be a neighbourhood  $\mathcal{U}$  of  $r$  such that  $\mathcal{U} \subset \mathcal{F}(q) \subset \mathcal{M} - \mathcal{N}$  which is impossible if  $r \in \dot{\mathcal{N}}$ . We may introduce a normal coordinate system  $x^a$  ( $a = 1, 2, 3, 4$ ) in a neighbourhood  $\mathcal{U}$  of  $q$  with  $x^4$  timelike, such that the lines  $x^i = \text{constant}$  ( $i = 1, 2, 3$ ) intersect  $\mathcal{P}(q)$  and  $\mathcal{F}(q)$ . Then each of these lines contains a point of  $\dot{\mathcal{N}}$ . Moreover, the  $x^4$  coordinates of these points must be continuous since no two points of  $\mathcal{N}$  have a timelike separation. Thus  $\dot{\mathcal{N}}$  is locally homeomorphic to  $E^3$ .

**LEMMA 5.** *If  $\mathcal{N}$  satisfies the conditions of lemma 4 and if for a point  $q \in \dot{\mathcal{N}}$  there is a normal coordinate neighbourhood  $\mathcal{U}$  of  $q$  and a ball  $\mathcal{B} \subset \mathcal{U}$  of constant coordinate radius about  $q$  which contains a sequence of points  $y_n$  converging to  $q$  such that there is a future directed null or timelike line  $\lambda_n$  from each point  $y_n$  which intersects  $(\mathcal{B} \cap \mathcal{N})$ , then  $\dot{\mathcal{N}}$  contains a future directed null geodesic segment from  $q$ .*

(The author is indebted to Dr R. Penrose for the following proof.)

Let  $z$  be a limit point of  $\lambda_n \cap (\mathcal{B} \cap \mathcal{N})$ . Select a subsequence  $y'_n$  such that  $\lambda'_n \cap (\mathcal{B} \cap \mathcal{N})$  converges to  $z$ . Then if  $\mathcal{V}$  is any neighbourhood of  $z$ , all  $y'_n$  for  $n$  large enough are contained in  $\overline{\mathcal{P}(\mathcal{V})}$ . Thus  $q \in \overline{\mathcal{P}(\mathcal{V})}$ . Therefore  $\mathcal{V}$  intersects  $\overline{\mathcal{F}(q)}$ . This shows that  $z \in \overline{\mathcal{F}(q)}$ . However  $z \in \dot{\mathcal{N}}$ , and thus it is not contained in  $\mathcal{F}(q)$ . Therefore  $z$  lies on the future null cone from  $q$  and the null geodesic segment  $qz$  lies in  $\dot{\mathcal{N}}$ .

**COROLLARY.** If there is a past directed null geodesic segment from  $q$  in  $\dot{\mathcal{N}}$  this must be a continuation of the future directed null geodesic segment from  $q$  since if it were in any other direction there would be points of  $\dot{\mathcal{N}}$  with timelike separation. If there is more than one future-directed null geodesic segment from  $q$  in  $\dot{\mathcal{N}}$  there can be no past directed null geodesic segment from  $q$  in  $\dot{\mathcal{N}}$ .

A boundary of a set satisfying lemmas 4 and 5 will be called a null horizon.

**LEMMA 6.** *For any  $p$  the set  $\mathcal{I}(p) \equiv \mathcal{F}(p) \cap \mathcal{P}(\mathcal{H})$  has compact closure.*

Let  $\mathcal{C}$  be the set of all points which have an open neighbourhood  $\mathcal{W}$  such that  $\overline{\mathcal{I}(\mathcal{W})}$  is compact or empty. Points near  $\mathcal{H}$  will be in  $\mathcal{C}$ . Suppose  $\mathcal{D} \equiv \mathcal{M} - \mathcal{C}$  is not empty. Then it satisfies the conditions of lemma 4. We shall see that the conditions of lemma 5 are satisfied at each point of  $\dot{\mathcal{D}}$ . Let  $q \in \dot{\mathcal{D}}$ , and let  $\mathcal{B}$  be a ball of constant coordinate radius about  $q$  in a normal coordinate neighbourhood  $\mathcal{U}$ . Let  $y_n$  be a sequence of points in  $\mathcal{P}(q)$  converging to  $q$  with  $y_{n+1} \in \mathcal{F}(y_n)$ . Let  $\mathcal{V}_n$  be a set of open neighbourhoods in  $\mathcal{P}(q)$  with  $y_{n+1} \in \mathcal{V}_{n+1} \subset \mathcal{F}(y_n)$ . Suppose that for some  $n$ ,  $(\overline{\mathcal{F}(\mathcal{V}_n)} \cap \mathcal{B}) \subset \mathcal{C}$  then since  $\overline{\mathcal{F}(\mathcal{V}_n)} \cap \mathcal{B}$  is compact it can be covered by a finite number of open neighbourhoods  $\mathcal{W}_i$  in  $\mathcal{C}$  for which  $\overline{\mathcal{I}(\mathcal{W}_i)}$  is compact. But

then  $\bar{\mathcal{I}}(\mathcal{V}_n)$  would be contained in  $[\bigcup_i \bar{\mathcal{I}}(\mathcal{W}_i)] \cup [\bar{\mathcal{F}}(\mathcal{V}_n) \cap \bar{\mathcal{B}}]$  which is compact.

This is impossible since it would imply  $y_n \in \mathcal{C}$ . Thus from each  $\mathcal{V}_n$  and each  $y_{n-1}$  there must be a timelike or null line which intersects  $\dot{\mathcal{B}} \cap \mathcal{D}$ . This shows, by lemma 5, that  $\dot{\mathcal{D}}$  is a null horizon, and that the null geodesic segments generating  $\mathcal{D}$  can have no future end-point. However,  $\dot{\mathcal{D}}$  will be in  $\mathcal{P}(\mathcal{H})$  and every future directed null geodesic in  $\mathcal{P}(\mathcal{H})$  intersects  $\mathcal{H}$  which is not in  $\mathcal{D}$ . This would imply that the generators of  $\dot{\mathcal{D}}$  had future end-points which shows that the supposition that  $\mathcal{D}$  was nonempty must be false. Thus  $\bar{\mathcal{I}}(p)$  is compact for all  $p$ .

For a point  $p \in \mathcal{P}(\mathcal{S})$  we define  $d(p, \mathcal{S})$  to be the least upper bound of the lengths of future directed timelike lines from  $p$  to  $\mathcal{S}$ . We define  $d(p, \mathcal{S})$  to be zero if  $p$  is not in  $\mathcal{P}(\mathcal{S})$ . Near the Cauchy surface  $\mathcal{H}$ , there will be points in  $\mathcal{P}(\mathcal{H})$  with  $d(p, \mathcal{H})$  finite.

**LEMMA 7.** *If  $p \in \mathcal{P}(\mathcal{H})$  and  $d(p, \mathcal{H})$  is finite, then there is a future directed timelike geodesic from  $p$  to  $\mathcal{H}$  of length  $d(p, \mathcal{H})$ .*

Let  $\mathcal{G}(a)$  be the set of all points  $r$  such that  $d(r, \mathcal{H}) = a$ . Since  $\bar{\mathcal{I}}(p)$  is compact we can find a sufficiently small  $a > 0$  such that every point  $q \in (\mathcal{G}(a) \cap \bar{\mathcal{I}}(p))$  is connected to  $\mathcal{H}$  by a unique future directed timelike geodesic of length  $a$ . Then  $\mathcal{G}(a) \cap \bar{\mathcal{I}}(p)$  is compact. Let  $q \in (\mathcal{G}(a) \cap \bar{\mathcal{I}}(p))$  be such that  $d(p, q)$  is maximized over  $\mathcal{G}(a) \cap \bar{\mathcal{I}}(p)$ . Then  $d(p, q) = d(p, \mathcal{H}) - a$ . Let  $\gamma(s)$  be the past directed timelike geodesic from  $\mathcal{H}$  to  $q$  with  $s = 0$  at  $\mathcal{H}$  and  $s = a$  at  $q$ . Let  $\mathcal{U} \subset \mathcal{P}(\mathcal{H})$  be a normal coordinate neighbourhood in which  $\gamma(s)$  lies on the  $x^4$  axis and let  $\mathcal{B} \subset \mathcal{U}$  be a ball of constant coordinate radius about  $q$ . Let  $y \equiv \gamma(s) \cap \mathcal{B}$  for  $s < a$  and  $y' \equiv \gamma(s) \cap \mathcal{B}$  for  $s > a$ . Let  $z \in \mathcal{B} \cap \bar{\mathcal{I}}(q)$  be such that  $d(p, z)$  is maximized. Then  $z$  must coincide with  $y'$  since if it did not  $d(z, y)$  would be greater than  $d(z, q) + d(q, y)$  and there would be a timelike line from  $\mathcal{H}$  through  $y$  and  $z$  to  $p$  of length greater than  $d(p, \mathcal{H})$ . Thus the relation  $d(p, \gamma(s)) = d(p, \mathcal{H}) - s$  holds along  $\gamma(s)$  from  $s = 0$  to  $z$ . But then it must hold for all values of  $s$  with  $0 \leq s \leq d(p, \mathcal{H})$  because if  $s_0$  were the least upper bound of the values of  $s$  for which it held we could apply the same construction at the point  $\gamma(s_0)$  and prove it held for  $s = s_0 + \delta$ . Thus the point  $\gamma(d(p, \mathcal{H}))$  either coincides with  $p$  or lies on  $\bar{\mathcal{F}}(p)$  not at  $p$ . An application of a similar construction shows the latter is impossible. Thus  $\gamma$  is a past directed timelike geodesic of length  $d(p, \mathcal{H})$  from  $\mathcal{H}$  to  $p$ . It will be orthogonal to  $\mathcal{H}$ .

**LEMMA 8.**  *$d(p, H)$  is finite.*

Let  $\mathcal{N}$  be the set of all points  $p$  for which  $d(p, \mathcal{H})$  is infinite. Suppose  $\mathcal{N}$  is non-empty. Then the conditions of lemmas 4 and 5 would be satisfied, and  $\mathcal{N}$  hence would be a null horizon. This is impossible since the generators of  $\mathcal{N}$  would intersect  $\mathcal{H}$ . Hence  $\mathcal{N}$  is empty.

We are now in a position to prove theorem 1. Conditions (1) and (4) imply that there is a point conjugate to  $\mathcal{H}$  on every past directed geodesic normal to  $\mathcal{H}$  within a distance  $3/c$ . But any point  $p \in \mathcal{P}(\mathcal{H})$  can be joined to  $\mathcal{H}$  by a future directed timelike geodesic orthogonal to  $\mathcal{H}$  and of length  $d(p, \mathcal{H})$ . By lemma 3 this geodesic cannot contain a point conjugate to  $\mathcal{H}$ . But any past-directed timelike geodesic can be extended to a length greater than  $3/c$ . Thus there are points  $p$  for which  $d(p, \mathcal{H}) > 3/c$ .

This is a contradiction which shows that the assumption that space-time is singularity free is incompatible with conditions (1) to (4).

This theorem depended on condition (4) which might not be observationally verifiable since if we were on the Cauchy surface ourselves we would not be able to investigate other parts of it to see whether the normals were indeed expanding. Thus a theorem that depended on a condition that could be checked by observation might seem preferable.

**THEOREM 2.** *Space-time cannot be singularity-free if conditions (1) to (3) are satisfied, and*

(5) *There is a point  $q$  on  $\mathcal{H}$  such that all past-directed timelike geodesics through  $q$  have a positive lower bound for convergence within an upper bounded distance from  $q$ . That is all past-directed geodesics through  $q$  start converging again in a finite distance.*

Condition (5) will be satisfied by a Robertson model. It will therefore be satisfied by a model that is sufficiently similar to a Robertson model at the present time. Since it only depends on the behaviour of geodesics inside the light cone of  $q$ , it could in principle be tested by observation.

Conditions (1) and (5) imply that all past-directed timelike geodesics through  $q$  have a point conjugate to  $q$  within an upper bound of distance from  $q$ . The proof is then similar to that for theorem 1.

Theorems 1 and 2 prove the occurrence of singularities in the past for expanding universe models. They could also be applied to contracting models to prove the occurrence of singularities in the future.

## 7. CAUCHY SURFACES

Theorems 1 and 2 depended on the existence of a Cauchy surface. Penrose has suggested (private communication) that the universe might not possess a global Cauchy surface. He instances the Kerr solution (cf. Carter 1966) as an example of a space-time where, although no spacelike surface intersects any timelike or null line more than once, it is impossible to find a spacelike surface that does intersect every line. Although it would mean modifying our ideas of determinism, we have no *a priori* reason for believing that the universe might not have similar properties. However, if our ideas of causality have any validity we ought to be able to find a complete connected spacelike surface  $\mathcal{H}$  which although not necessarily intersecting every timelike and null line does not intersect any more than once. We shall call such a surface a partial Cauchy surface.  $\mathcal{K}$ , the set of points  $p$  such that every future directed timelike or null line from  $p$  intersects  $\mathcal{H}$  will be called the past Cauchy development of  $\mathcal{H}$ . If  $\mathcal{P}(\mathcal{H}) - \mathcal{H}$  is not empty we can apply lemmas 4 and 5 to show that  $\mathcal{K} - \mathcal{H}$  is a null horizon which we shall call the past Cauchy horizon of  $\mathcal{H}$ . This does not imply any contradiction since the partial Cauchy surface,  $\mathcal{H}$ , does not intersect every null line.

**THEOREM 3.** *Space-time cannot be singularity free if conditions (1) and (2) are satisfied, and*

(6) *There exists a smooth compact partial Cauchy surface  $\mathcal{H}$  whose unit normal  $v^a$  is expanding everywhere on  $\mathcal{H}$ , i.e.  $v^a_{;a} > 0$  on  $\mathcal{H}$ .*

(7) *There is an open neighbourhood  $\mathcal{W}$  of each point  $p \in \mathcal{P}(\mathcal{H})$  such that every future-directed null geodesic which intersects  $\mathcal{W}$  leaves it at some point and never reenters it. (In effect this is a condition that there should not be closed null lines in  $\mathcal{P}(\mathcal{H})$ .)*

Note that it is necessary that  $\mathcal{H}$  be compact since Minkowski space satisfies conditions (1), (2) and (7) and has a non-compact partial Cauchy surface whose normal is expanding.

LEMMA 9. *There exists a global nonvanishing past directed timelike  $C^2$  vector field  $u^a$  on  $\mathcal{M}$ . (The author is indebted to Dr R. Penrose and Professor C. W. Misner for this.)*

By a standard theorem (cf. Bishop & Crittenden 1964, p. 126),  $\mathcal{M}$  being paracompact may be endowed with a positive definite differentiable Riemannian metric  $m_{ab}$ . Consider the eigenvalue equation,  $m_{ab}u^b = \lambda g_{ab}u^b$  at a point  $p$ . By transforming to coordinates in which  $m_{ab}$  is the unit matrix at  $p$ , it can be seen that the timelike eigenvector  $u^a$  is unique since it corresponds to a positive eigenvalue, whereas the other eigenvalues are negative. The eigenvector  $u^a$  may be normalized by  $u^a u^b g_{ab} = 1$ .

This vector field may be used to give a diffeomorphism  $\beta$  of  $\mathcal{H} \times [0, t]$  into  $\mathcal{M}$ , for any  $t > 0$ ,\* by taking points of  $\mathcal{H}$  a distance  $s \in [0, t]$  along the past directed integral curves of  $u^a$ . If  $d(r, \mathcal{H})$  had an upper bound  $t$  for  $r \in \mathcal{K} - \mathcal{H}$  then

$$\mathcal{K} - \mathcal{H} \subset \beta(\mathcal{H} \times [0, t])$$

since every future directed timelike line through  $\mathcal{K} - \mathcal{H}$  intersects  $\mathcal{H}$ . Thus  $\mathcal{K} - \mathcal{H}$  would be compact and could be covered by a finite number of the neighbourhoods  $\mathcal{W}_i$  of condition (7). A future directed null geodesic segment generating  $\mathcal{K} - \mathcal{H}$  must leave each  $\mathcal{W}_i$  it enters and not reenter. This is impossible since it cannot have a future end-point. Hence  $d(r, \mathcal{H})$  can have no upper bound for  $r \in \mathcal{K} - \mathcal{H}$ .

By conditions (1) and (6) there is a point conjugate to  $\mathcal{H}$  on each past directed geodesic normal to  $\mathcal{H}$  within some finite upper bound  $b$  of distance from  $\mathcal{H}$ . Let  $\mathcal{L}$  be the set of all points  $p$  for which  $d(p, \mathcal{H}) > b$ . Then  $\mathcal{L} \cap (\mathcal{K} - \mathcal{H})$  is not empty. Since  $\mathcal{L}$  is open,  $\mathcal{L} \cap \mathcal{H}$  is not empty. If  $r \in \mathcal{K}$  lemmas 6, 7 and 8 may be applied to show that  $d(r, \mathcal{H})$  is finite and that  $r$  may be joined to  $\mathcal{H}$  by a future directed line of length  $d(r, \mathcal{H})$ . This leads to the same contradiction as in theorem 1 since  $p \in \mathcal{L} \cap \mathcal{H}$  would be joined to  $\mathcal{H}$  by a line of length  $d(p, \mathcal{H}) > b$ .

## 8. CONCLUSION

The preceding sections indicate that singularities are inevitable in certain models if condition (1) is satisfied and the Einstein equations hold. Unlike the singularities in collapsed stars which cannot be seen by an external observer, these singularities would in principle be observable. So far the Einstein equations have only been experimentally tested for very weak fields for which they give almost the same results as Newtonian theory. Thus observations to determine whether singularities actually occurred would provide a test of the equations for strong fields. Presumably

\* *Note added in proof.* This is incorrect. A correct proof is possible using the diffeomorphism defined by the geodesics normal to  $\mathcal{H}$ .

it would be necessary to consider quantum effects in very strong fields. However, these would not become important until the radius of curvature became of the order of  $10^{-14}$  cm which for practical purposes is pretty singular.

The view has been expressed that singularities are so objectionable that if the Einstein equations were to predict their occurrence, this would be a compulsive reason for modifying them. However, the real test of a physical theory is not whether its predicted results are aesthetically attractive but whether they agree with observation. So far there are no observations which would show that singularities do not occur.

Conditions (1) will be satisfied by normal matter. However, it will not be satisfied by the 'C' field of Hoyle & Narlikar (1963). This is a field of negative energy density which can prevent the occurrence of singularities. However, there would seem to be a certain quantum-mechanical difficulty associated with the existence of fields of negative energy density. For there would not seem anything to prevent the creation, in a given volume of space-time, of an infinite number of particles of positive energy and a corresponding infinite number of quanta of the negative energy field. It might be possible to overcome this difficulty but so far no workable scheme has been proposed.

Theorems 1, 2 and 3 prove the existence of a singularity but give no information as to its nature. Thus it need not be like the all-embracing singularity in the Robertson models which every world-line hits: it might be an isolated singularity which only a few lines hit. Further research will be needed to determine the actual character of the singularity in each situation.

The author is grateful to Mr B. Carter, Dr G. R. F. Ellis, Dr R. Penrose and Dr D. W. Sciama for many useful discussions and helpful suggestions.

#### REFERENCES

- Bishop, R. L. & Crittenden, R. J. 1964 *The geometry of manifolds*. Academic Press.  
 Carter, B. 1966 *Phys. Rev.* **141**, 1242.  
 Hawking, S. W. 1965 *Phys. Rev. Lett.* **15**, 689.  
 Hawking, S. W. & Ellis, G. R. F. 1965 *Phys. Lett.* **17**, 246.  
 Hoyle, F. & Narlikar, J. V. 1963 *Proc. Roy. Soc. A* **278**, 465.  
 Komar, A. 1956 *Phys. Rev.* **104**, 544.  
 Lifshitz, E. M. & Khalatnikov, I. M. 1963 *Adv. Phys.* **12**, 183.  
 Milnor, J. 1963 *Morse theory*. Princeton University Press.  
 Penrose, R. 1965 *Phys. Rev. Lett.* **14**, 57.  
 Raychaudhuri, A. 1955 *Phys. Rev.* **98**, 1123.  
 Robertson, H. P. 1933 *Rev. Mod. Phys.* **5**, 62.  
 Shepley, L. C. 1964 *Proc. Nat. Acad. Sci. U.S.A.* **52**, 1403.