

Notes on Lorentzian causality

ESI-EMS-IAMP Summer School on Mathematical Relativity

Gregory J. Galloway
Department of Mathematics
University of Miami

August 4, 2014

Contents

1	Lorentzian manifolds	2
2	Futures and pasts	8
3	Achronal boundaries	11
4	Causality conditions	14
5	Domains of dependence	20
6	The geometry of null hypersurfaces	23
7	Trapped surfaces and the Penrose Singularity Theorem	29

1 Lorentzian manifolds

In General Relativity, the space of events is represented by a *Lorentzian manifold*, which is a smooth manifold M^{n+1} equipped with a metric g of Lorentzian signature. Thus, at each $p \in M$,

$$g : T_p M \times T_p M \rightarrow \mathbb{R} \quad (1.1)$$

is a scalar product of signature $(-, +, \dots, +)$. With respect to an orthonormal basis $\{e_0, e_1, \dots, e_n\}$, as a matrix,

$$[g(e_i, e_j)] = \text{diag}(-1, +1, \dots, +1). \quad (1.2)$$

Example: Minkowski space, the spacetime of Special Relativity. Minkowski space is \mathbb{R}^{n+1} , equipped with the Minkowski metric η : For vectors $X = X^i \frac{\partial}{\partial x^i}$, $Y = Y^i \frac{\partial}{\partial x^i}$ at p , (where x^i are standard Cartesian coordinates on \mathbb{R}^{n+1}),

$$\eta(X, Y) = -X^0 Y^0 + \sum_{i=1}^n X^i Y^i. \quad (1.3)$$

Similarly, for the Lorentzian metric g , we have for vectors $X = X^i e_i$, $Y = Y^j e_j$ at p ,

$$g(X, Y) = g(e_i, e_j) X^i Y^j = -X^0 Y^0 + \sum_{i=1}^n X^i Y^i. \quad (1.4)$$

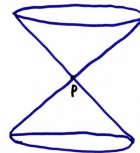
Thus, each tangent space of a Lorentzian manifold is isometric to Minkowski space. Hence, one may say that Lorentzian manifolds are locally modeled on Minkowski space, just as Riemannian manifolds are locally modeled on Euclidean space.

1.1 Causal character of vectors.

At each point, vectors fall into three classes, as follows:

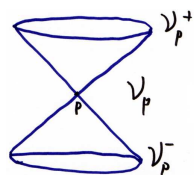
$$X \text{ is } \begin{cases} \text{timelike} & \text{if } g(X, X) < 0 \\ \text{null} & \text{if } g(X, X) = 0 \\ \text{spacelike} & \text{if } g(X, X) > 0. \end{cases}$$

We see that the set of null vectors $X \in T_p M$ forms a double cone \mathcal{V}_p in the tangent space $T_p M$:

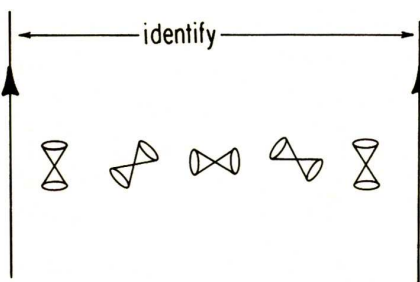


called the null cone (or light cone) at p . Timelike vectors point inside the null cone and spacelike vectors point outside.

Time orientability. Consider at each point of p in a Lorentzian manifold M the null cone $\mathcal{V}_p \subset T_pM$. \mathcal{V}_p is a double cone consisting of two cones, \mathcal{V}_p^+ and \mathcal{V}_p^- :



We may designate one of the cones, \mathcal{V}_p^+ , say, as the *future* null cone at p , and the other half cone, \mathcal{V}_p^- , as the *past* null cone at p . If this assignment can be made in a continuous manner over all of M (this can always be done locally) then we say that M is *time-orientable*. The following figure illustrates a Lorentzian manifold that is *not* time-orientable (even though the underlying manifold is orientable).



There are various ways to make the phrase “continuous assignment” precise (see e.g., [15, p. 145]), but they all result in the following, which we adopt as the definition of time-orientability.

Definition 1.1. *A Lorentzian manifold M^{n+1} is time-orientable iff it admits a smooth timelike vector field T .*

If M is time-orientable, the choice of a smooth time-like vector field T fixes a time orientation on M . For any $p \in M$, a (nonzero) causal vector $X \in T_pM$ is future directed (resp. past directed) provided $g(X, T) < 0$ (resp. $g(X, T) > 0$). Thus X is future directed if it points into the same half cone at p as T .

We remark that if M is not time-orientable, it admits a double cover that is.

By a *spacetime* we mean a connected time-oriented Lorentzian manifold (M^{n+1}, g) . We will usually restrict attention to spacetimes.

Lorentzian inequalities. We say that $X \in T_pM$ is *causal* if it is time like or null, $g(X, X) \leq 0$. If X is causal, define its length as

$$|X| = \sqrt{-g(X, X)}.$$

Proposition 1.1. *The following basic inequalities hold.*

(1) (Reverse Schwarz inequality) For all causal vectors $X, Y \in T_p M$,

$$|g(X, Y)| \geq |X||Y| \quad (1.5)$$

(2) (Reverse triangle inequality) For all causal vectors X, Y that point into the same half cone of the null cone at p ,

$$|X + Y| \geq |X| + |Y|. \quad (1.6)$$

Proof hints: Note (1.5) trivially holds if X is null. For X timelike, decompose Y as $Y = \lambda X + Y^\perp$, where Y^\perp (the component of Y perpendicular to X) is necessarily spacelike. Inequality (1.6) follows easily from (1.5).

The Reverse triangle inequality is the source of the twin paradox.

1.2 Causal character of curves:

Let $\gamma : I \rightarrow M, t \rightarrow \gamma(t)$ be a smooth curve in M .

γ is said to be *timelike* provided $\gamma'(t)$ is timelike for all $t \in I$.



In GR, a timelike curve corresponds to the history (or *worldline*) of an observer.

Null curves and *spacelike curves* are defined analogously.

A *causal curve* is a curve whose tangent is either timelike or null at each point. Heuristically, in accordance with relativity, information flows along causal curves, and so such curves are the focus of our attention.

The notion of a causal curve extends in a natural way to piecewise smooth curves. The only extra requirement is that when two segments join, at some point p , say, the end point tangent vectors must point into the same half cone of the null cone \mathcal{V}_p at p . We will normally work within this class of piecewise smooth causal curves.

The length of a causal curve $\gamma : [a, b] \rightarrow M$, is defined by

$$L(\gamma) = \text{Length of } \gamma = \int_a^b |\gamma'(t)| dt = \int_a^b \sqrt{-\langle \gamma'(t), \gamma'(t) \rangle} dt.$$

If γ is timelike one can introduce arc length parameter along γ . In general relativity, the arc length parameter along a timelike curve is called proper time, and corresponds to time kept by the observer.

1.3 The Levi-Civita connection and geodesics.

Recall that a Lorentzian manifold M (like any pseudo-Riemannian manifold) admits a unique covariant derivative operator ∇ called the *Levi-Civita connection*. Thus for smooth vector fields X, Y on M , $\nabla_X Y$ is a vector field on M (the directional derivative of Y in the direction X) satisfying:

- (1) $\nabla_X Y$ is linear in Y over the reals.
- (2) $\nabla_X Y$ is linear in Y over the space of smooth functions. (In particular, $\nabla_{fX} Y = f\nabla_X Y$).
- (3) (Product rule) $\nabla_X fY = X(f)Y + f\nabla_X Y$.
- (4) (Symmetric) $[X, Y] = \nabla_X Y - \nabla_Y X$.
- (5) (Metric product rule) $X\langle Y, Z \rangle = \langle \nabla_X Y, Z \rangle + \langle Y, \nabla_X Z \rangle$.

∇ is uniquely determined by these properties. With respect to a coordinate neighborhood (U, x^i) , one has,

$$\nabla_X Y = (X(Y^k) + \Gamma_{ij}^k X^i Y^j) \partial_k, \quad (1.7)$$

where $\partial_i = \frac{\partial}{\partial x^i}$, $X = X^i \partial_i$, $Y = Y^j \partial_j$, and the Γ_{ij}^k 's are the Christoffel symbols,

$$\Gamma_{ij}^k = \frac{1}{2} g^{km} (g_{jm,i} + g_{im,j} - g_{ij,m}),$$

where $g_{ij} = g(\partial_i, \partial_j)$, etc.

We see from the coordinate expression in (1.7) that $\nabla_X Y$ depends only on the value of X at a point and only on the values of Y along a curve, defined in neighborhood of the point, having X as a tangent vector.

Thus the Levi-Civita connection enables one to compute the covariant derivative of a vector field $t \xrightarrow{Y} Y(t) \in T_{\gamma(t)} M$ defined along a curve $\gamma : I \rightarrow M$, $t \rightarrow \gamma(t)$. In local coordinates $\gamma(t) = (x^1(t), \dots, x^n(t))$, and from (1.7) we have

$$\nabla_{\gamma'} Y = \left(\frac{dY^k}{dt} + \Gamma_{ij}^k \frac{dx^i}{dt} Y^j \right) \partial_k. \quad (1.8)$$

where $\gamma' = \frac{dx^i}{dt} \partial_i|_{\gamma}$ is the tangent (or velocity) vector field along γ and $Y(t) = Y^i(t) \partial_i|_{\gamma(t)}$.

Geodesics. Given a curve $t \rightarrow \gamma(t)$ in M , $\nabla_{\gamma'} \gamma'$ is called the *covariant acceleration* of γ . In local coordinates,

$$\nabla_{\gamma'} \gamma' = \left(\frac{d^2 x^k}{dt^2} + \Gamma_{ij}^k \frac{dx^i}{dt} \frac{dx^j}{dt} \right) \partial_k, \quad (1.9)$$

as follows by setting $Y^k = \frac{dx^k}{dt}$ in Equation (1.8). By definition, a *geodesic* is a curve of zero covariant acceleration,

$$\nabla_{\gamma'} \gamma' = 0 \quad (\text{Geodesic equation}) \quad (1.10)$$

In local coordinates the geodesic equation becomes a system of $n + 1$ second order ODE's in the coordinate functions $x^i = x^i(t)$,

$$\frac{d^2 x^k}{dt^2} + \Gamma_{ij}^k \frac{dx^i}{dt} \frac{dx^j}{dt} = 0, \quad k = 0, \dots, n. \quad (1.11)$$

The basic existence and uniqueness result for systems of ODE's guarantees the following.

Proposition 1.2. *Given $p \in M$ and $v \in T_p M$, there exists an interval I about $t = 0$ and a unique geodesic $\sigma : I \rightarrow M$, $t \rightarrow \sigma(t)$, satisfying,*

$$\sigma(0) = p, \quad \frac{d\sigma}{dt}(0) = v.$$

In fact, by a more refined analysis it can be shown that each $p \in M$ is contained in a (*geodesically*) *convex* neighborhood U , which has the property that any two points in U can be joined by a unique geodesic contained in U . In fact U can be chosen so as to be a normal neighborhood of each of its points; cf. [15], p. 129. (Recall, a normal neighborhood of $p \in M$ is the diffeomorphic image under the exponential map of a star-shaped domain about $\mathbf{0} \in T_p M$.)

Finally, note if γ is a geodesic then by the metric product rule, $\gamma'(g(\gamma', \gamma')) = 2g(\nabla_{\gamma'} \gamma', \gamma') = 0$, and hence geodesics are always constant speed curves. Thus, each geodesic in a Lorentzian manifold is either timelike, spacelike or null. In GR timelike geodesics correspond to *freely falling* observers and null geodesics correspond to the paths of photons.

1.4 Local Lorentz geometry. In Minkowski space the geodesics are straight lines (the Christoffel symbols vanish in Cartesian coordinates). Moreover the following holds:

- (1) If there is a timelike curve γ from p to q then \overline{pq} is timelike.
- (2) $L(\overline{pq}) \geq L(\gamma)$, for all causal curves γ from p to q .

Although it can be very different in the large, *locally* the geometry and causality of a Lorentzian manifold is similar to Minkowski space. Let U be a convex neighborhood in a Lorentzian manifold. Hence for each pair of points $p, q \in U$ there exists a unique geodesic segment from p to q in U , which we denote by \overline{pq} .

Proposition 1.3 ([15], p. 146). *Let U be a convex neighborhood in a Lorentzian manifold M^{n+1} .*

- (1) *If there is a timelike (resp., causal) curve in U from p to q then \overline{pq} is timelike (causal).*
- (2) *If \overline{pq} is timelike then $L(\overline{pq}) \geq L(\gamma)$ for all causal curves γ in U from p to q . Moreover, the inequality is strict unless, when suitable parametrized, $\gamma = \overline{pq}$.*

Thus, within a convex neighborhood U , timelike geodesics are *maximal*, i.e., are causal curves of greatest length. Moreover, within U null geodesics are *achronal*, i.e., no two points can be joined by a timlike curve. Both of these features can fail in the large.

1.5 Curvature and the Einstein equations

The Riemann curvature tensor of (M, g) is defined in terms of second covariant derivatives anti-symmetrized: For vector fields $X, Y, Z \in \mathfrak{X}(M)$,

$$R(X, Y)Z = \nabla_X \nabla_Y Z - \nabla_Y \nabla_X Z - \nabla_{[X, Y]} Z. \quad (1.12)$$

The components $R^l{}_{kij}$ of the Riemann curvature tensor R in a coordinate chart (U, x^i) are defined by the following equation,

$$R(\partial_i, \partial_j)\partial_k = R^l{}_{kij}\partial_l$$

The Ricci tensor is obtained by contraction,

$$R_{ij} = R^l{}_{ilj}$$

Symmetries of the Riemann curvature tensor imply that the Ricci tensor is symmetric, $R_{ij} = R_{ji}$. By tracing the Ricci tensor, we obtain the scalar curvature,

$$R = g^{ij} R_{ij}.$$

The Einstein equations, the field equations of GR, are given by:

$$R_{ij} - \frac{1}{2} R g_{ij} = 8\pi T_{ij},$$

where T_{ij} is the energy-momentum tensor. The Einstein equations describe how spacetime curves in the presence of matter, and it is this curvature that is responsible for the effects of gravity. The left hand side is a purely geometric tensor, the Einstein tensor. On the right hand side is the energy momentum tensor T , which describes the energy-momentum content of matter and all other nongravitational fields.

From the PDE point of view, the Einstein equations form a system of second order quasi-linear equations for the g_{ij} 's. This system may be viewed as a (highly complicated!) generalization of Poisson's equation in Newtonian gravity.

The *vacuum* Einstein equations are obtained by setting $T_{ij} = 0$. It is easily seen that this equivalent to setting $R_{ij} = 0$. We will sometimes require that a space-time satisfying the Einstein equations, obeys an *energy condition*. The *null energy condition* (NEC) is the requirement that

$$T(X, X) = \sum_{i,j} T_{ij} X^i X^j \geq 0 \quad \text{for all null vectors } X. \quad (1.13)$$

The stronger *dominant energy condition* (DEC) is the requirement,

$$T(X, Y) = \sum_{i,j} T_{ij} X^i Y^j \geq 0 \quad \text{for all future directed causal vectors } X, Y. \quad (1.14)$$

The DEC is satisfied by most classical fields. Physically, the DEC requires that the speed of energy flow is less than the speed of light.

2 Futures and pasts

We begin the study of causal theory in earnest. Causal theory is the study of the causal relations ' \ll ' and ' $<$ '.

Let (M, g) be a spacetime. A timelike (resp. causal) curve $\gamma : I \rightarrow M$ is said to be *future directed* provided each tangent vector $\gamma'(t)$, $t \in I$, is future directed. (*Past-directed* timelike and causal curves are defined in a time-dual manner.)

Definition 2.1. For $p, q \in M$,

- (1) $p \ll q$ means there exists a future directed timelike curve in M from p to q (we say that q is in the timelike future of p),
- (2) $p < q$ means there exists a future directed (nontrivial) causal curve in M from p to q (we say that q is in the causal future of p),

We shall use the notation $p \leq q$ to mean $p = q$ or $p < q$.

The causal relations \ll and $<$ are clearly transitive. Also, from variational considerations, it is heuristically clear that the following holds,

$$\text{if } p \ll q \text{ and } q < r \text{ or if } p < q \text{ and } q \ll r \text{ then } p \ll r.$$



The above statement is a consequence of the following fundamental causality result; see [15, p. 294] for a careful proof.

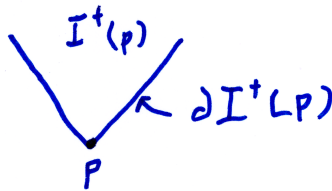
Proposition 2.1. In a spacetime M , if q is in the causal future of p ($p < q$) but is not in the timelike future of p ($p \not\ll q$) then any future directed causal curve γ from p to q must be a null geodesic (when suitably parameterized).

Definition 2.2. Given any point p in a spacetime M , the timelike future and causal future of p , denoted $I^+(p)$ and $J^+(p)$, respectively, are defined as,

$$I^+(p) = \{q \in M : p \ll q\} \quad \text{and} \quad J^+(p) = \{q \in M : p \leq q\}.$$

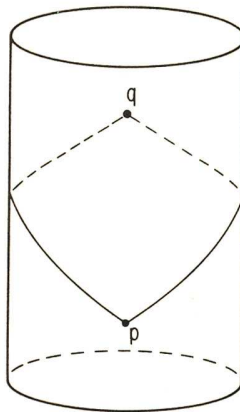
Hence, $I^+(p)$ consists of all points in M that can be reached from p by a future directed timelike curve, and $J^+(p)$ consists of the point p and all points in M that can be reached from p by a future directed causal curve. The timelike and causal *pasts* of p , $I^-(p)$ and $J^-(p)$, respectively, are defined in a time-dual manner in terms of past directed timelike and causal curves. Note by Proposition 2.1, if $q \in J^+(p) \setminus I^+(p)$ ($q \neq p$) then there exists a future directed null geodesic from p to q .

Ex. Minkowski space. For p any point in Minkowski space, $I^+(p)$ is open, $J^+(p)$ is closed and $\partial I^+(p) = J^+(p) \setminus I^+(p)$ is just the future null cone at p . $I^+(p)$ consists of all points inside the future null cone, and $J^+(p)$ consists of all points on and inside the future null cone.



We note, however, that curvature and topology can drastically change the structure of ‘null cones’ in spacetime.

Ex. Consider the example depicted in the following figure of a flat spacetime cylinder, closed in space. For any point p in this spacetime the future ‘null cone’ at p , $\partial I^+(p)$, is compact and consists of the two future directed null geodesic segments emanating from p that meet to the future at a point q . By extending these geodesics beyond q we enter $I^+(p)$.



In some situations it is convenient to restrict the causal relations \ll and $<$ to open subsets U of a spacetime M . For example, $I^+(p, U)$, the chronological future of p within U , consists of all points q in U for which there exists a future directed timelike curve *within* U from p to q , etc. Note that, in general $I^+(p, U) \neq I^+(p) \cap U$.

In general the sets $I^+(p)$ in a spacetime M are open. This is heuristically rather clear: A sufficiently small smooth perturbation of a timelike curve is still timelike. A rigorous proof is based on the causality of convex neighborhoods.

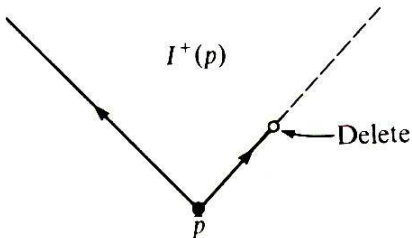
Proposition 2.2. *Let U be a convex neighborhood in a spacetime M . Then, for each $p \in U$,*

- (1) $I^+(p, U)$ is open in U (and hence M),
- (2) $J^+(p, U)$ is the closure in U of $I^+(p, U)$.

This proposition follows essentially from part (1) of Proposition 1.3.

Exercise: Prove that for each p in a spacetime M , $I^+(p)$ is open.

In general, sets of the form $J^+(p)$ need not be closed. This can be seen by removing a point from Minkowski space, as illustrated in the figure below.



Points on the dashed line are not in $J^+(p)$, but are in the closure $\overline{J^+(p)}$.

For any subset $S \subset M$, we define the timelike and causal future of S , $I^+(S)$ and $J^+(S)$, respectively by

$$I^+(S) = \bigcup_{p \in S} I^+(p) \quad \text{and} \quad J^+(S) = \bigcup_{p \in S} J^+(p).$$

Thus, $I^+(S)$ consists of all points in M reached by a future directed timelike curve starting from S , and $J^+(S)$ consists of the points of S , together with the points in M reached by a future directed causal curve starting from S . Since arbitrary unions of open sets are open, it follows that $I^+(S)$ is always an open set. $I^-(S)$ and $J^-(S)$ are defined in a time-dual manner.

Although in general $J^+(S) \neq \overline{I^+(S)}$, the following relationships always hold between $I^+(S)$ and $J^+(S)$.

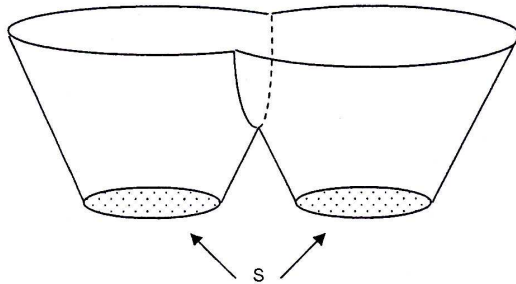
Proposition 2.3. *For all subsets $S \subset M$,*

- (1) $\text{int } J^+(S) = I^+(S)$,
- (2) $J^+(S) \subset \overline{I^+(S)}$.

Proof. Exercise.

3 Achronal boundaries

Achronal sets play an important role in causal theory. A subset $S \subset M$ is *achronal* provided no two of its points can be joined by a timelike curve. Of particular importance are *achronal boundaries*. By definition, an achronal boundary is a set of the form $\partial I^+(S)$ (or $\partial I^-(S)$), for some $S \subset M$. We wish to describe several important structural properties of achronal boundaries. The following figure illustrates nicely the properties to be discussed. It depicts the achronal boundary $\partial I^+(S)$ in Minkowski 3-space \mathbb{M}^3 , where S is the disjoint union of two spacelike disks; $\partial I^+(S)$ consists of S and the merging of two future light cones.



Proposition 3.1. *An achronal boundary $\partial I^+(S)$, if nonempty, is a closed achronal C^0 hypersurface in M .*

We discuss the proof of this proposition, beginning with the following basic lemma.

Lemma 3.2. *If $p \in \partial I^+(S)$ then $I^+(p) \subset I^+(S)$, and $I^-(p) \subset M \setminus \overline{I^+(S)}$.*

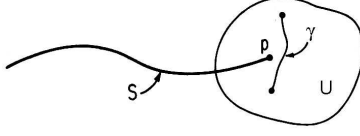
Proof. To prove the first part of the lemma, note that if $q \in I^+(p)$ then $p \in I^-(q)$, and hence $I^-(q)$ is a neighborhood of p . Since p is on the boundary of $I^+(S)$, it follows that $I^-(q) \cap I^+(S) \neq \emptyset$, and hence $q \in I^+(S)$. The second part of the lemma, which can be proved similarly, is left as an exercise. \square

Claim 1: An achronal boundary $\partial I^+(S)$ is achronal.

Proof of the claim: Suppose there exist $p, q \in \partial I^+(S)$, with $q \in I^+(p)$. By Lemma 3.2, $q \in I^+(S)$. But since $I^+(S)$ is open, $I^+(S) \cap \partial I^+(S) = \emptyset$. Thus, $\partial I^+(S)$ is achronal. \square

Lemma 3.2 also implies that achronal boundaries are *edgeless*. We need to introduce the edge concept.

Definition 3.1. *Let $S \subset M$ be achronal. Then $p \in \overline{S}$ is an edge point of S provided every neighborhood U of p contains a timelike curve γ from $I^-(p, U)$ to $I^+(p, U)$ that does not meet S (see the figure).*



We denote by $\text{edge } S$ the set of edge points of S . Note that $\overline{S} \setminus S \subset \text{edge } S \subset \overline{S}$. If $\text{edge } S = \emptyset$ we say that S is edgeless.

Claim 2: An achronal boundary $\partial I^+(S)$ is edgeless.

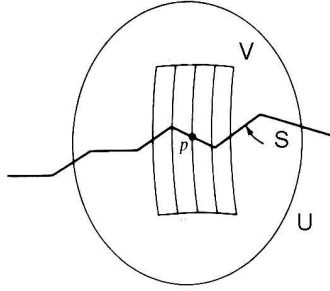
Proof of the claim: Lemma 3.2 implies that for any $p \in \partial I^+(S)$, any timelike curve from $I^-(p)$ to $I^+(p)$ must meet $\partial I^+(S)$. It follows that $\partial I^+(S)$ is edgeless. \square

Proposition 3.1 now follows from the following basic result.

Proposition 3.3. *Let S be achronal. Then $S \setminus \text{edge } S$, if nonempty, is a C^0 hypersurface in M . In particular, an edgeless achronal set is a C^0 hypersurface in M .*

Proof. We sketch the proof; for details, see [15, p. 413]. It suffices to show that in a neighborhood of each $p \in S \setminus \text{edge } S$, $S \setminus \text{edge } S$ can be expressed as a C^0 graph over a smooth hypersurface.

Fix $p \in S \setminus \text{edge } S$. Since p is not an edge point there exists a neighborhood U of p such that every timelike curve from $I^-(p, U)$ to $I^+(p, U)$ meets S . Let X be a future directed timelike vector field on M , and let \mathcal{N} be a smooth hypersurface in U transverse to X near p . Then, by choosing \mathcal{N} small enough, each integral curve of X through \mathcal{N} will meet S , and meet it exactly once since S is achronal. Using the flow generated by X , it follows that there is a neighborhood $V \approx (t_1, t_2) \times \mathcal{N}$ of p such that $S \cap V$ is given as the graph of a function $t = h(x)$, $x \in \mathcal{N}$ (see the figure below)



One can now show that a discontinuity of h at some point $x_0 \in \mathcal{N}$ leads to an achronality violation of S . Hence h must be continuous. \square

The next result shows that, in general, large portions of achronal boundaries are ruled by null geodesics. A future (resp., past) directed causal curve $\gamma : (a, b) \rightarrow M$ is said to be *future (resp., past) inextendible* in M if $\lim_{t \rightarrow b^-} \gamma(t)$ does not exist. A future directed causal curve $\gamma : (a, b) \rightarrow M$ is said to be inextendible if γ and $-\gamma$ are future and past inextendible, respectively.

Proposition 3.4. *Let $S \subset M$ be closed. Then each $p \in \partial I^+(S) \setminus S$ lies on a null geodesic contained in $\partial I^+(S)$, which either has a past end point on S , or else is past inextendible in M .*

The proof uses a standard tool in causal theory, namely that of taking a limit of causal curves. A technical difficulty arises however in that a limit of smooth causal curves need not be smooth. Thus, we are led to introduce the notion of a C^0 causal curve.

Definition 3.2. *A continuous curve $\gamma : I \rightarrow M$ is said to be a future directed C^0 causal curve provided for each $t_0 \in I$, there is an open subinterval $I_0 \subset I$ about t_0 and a convex neighborhood U of $\gamma(t_0)$ such that given any $t_1, t_2 \in I_0$ with $t_1 < t_2$, there exists a smooth future directed causal curve in U from $\gamma(t_1)$ to $\gamma(t_2)$.*

Thus, a C^0 causal curve is a continuous curve that can be approximated with arbitrary precision by a piecewise smooth causal curve.

We now give a version of the *limit curve lemma* (cf., [1, p. 511]). For its statement it is convenient to introduce a background complete Riemannian (positive definite) metric h on M . Observe that any future inextendible causal γ will have infinite length to the future, as measured in the metric h . Hence, if parameterized with respect to h -arc length, γ will be defined on the interval $[0, \infty)$.

Lemma 3.5 (Limit curve lemma). *Let $\gamma_n : [0, \infty) \rightarrow M$ be a sequence of future inextendible causal curves, parameterized with respect to h -arc length, and suppose that $p \in M$ is an accumulation point of the sequence $\{\gamma_n(0)\}$. Then there exists a future inextendible C^0 causal curve $\gamma : [0, \infty) \rightarrow M$ such that $\gamma(0) = p$ and a subsequence $\{\gamma_m\}$ which converges to γ uniformly with respect to h on compact subsets of $[0, \infty)$.*

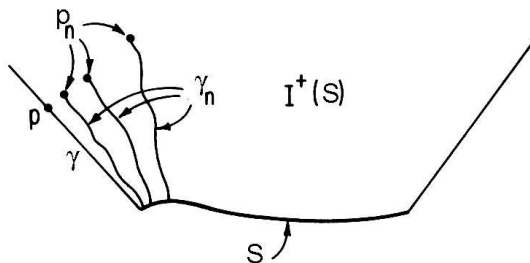
The proof of this lemma is an application of Arzela's theorem; see especially the proof of Proposition 3.31 in [1]. There are analogous versions of the limit curve lemma for past inextendible, and (past and future) inextendible causal curves.

Remark: We note that C^0 causal curves obey a local Lipschitz condition, and hence are rectifiable. Thus, in the limit curve lemma, the γ_n 's could be taken to be C^0 causal curves.¹ We also note that the "limit" parameter acquired by the limit curve γ need not in general be the h -arc length parameter.

Proof of Proposition 3.4. Fix $p \in \partial I^+(S) \setminus S$. Since $p \in \partial I^+(S)$, there exists a sequence of points $p_n \in I^+(S)$, such that $p_n \rightarrow p$. For each n , let $\gamma_n : [0, a_n] \rightarrow M$ be a past directed timelike curve from p_n to $q_n \in S$, parameterized with respect to h -arc length. Extend each γ_n to a past inextendible timelike curve $\tilde{\gamma}_n : [0, \infty) \rightarrow M$, parameterized with respect to h -arc length. By the limit curve lemma, there exists a subsequence $\tilde{\gamma}_m : [0, \infty) \rightarrow M$ that converges to a past inextendible C^0 causal curve $\gamma : [0, \infty) \rightarrow M$ such that $\gamma(0) = p$. By taking a further subsequence if necessary

¹See [5] for a treatment of causal theory based entirely on Lipschitz curves.

we can assume $a_m \uparrow a$, $a \in (0, \infty]$. We claim that $\gamma|_{[0,a]}$ (or $\gamma|_{[0,a)}$ if $a = \infty$) is the desired null geodesic (see the figure).



Fix $t \in (0, a)$. Eventually $a_m > t$, and so for large m we have $\tilde{\gamma}_m(t) = \gamma_m(t) \in I^+(S)$. Hence, since $\gamma(t) = \lim_{m \rightarrow \infty} \gamma_m(t)$, it follows that $\gamma(t) \in \overline{I^+(S)}$. Suppose $\gamma(t) \in I^+(S)$. Then there exists $x \in S$ such that $x \ll \gamma(t) < p$. This implies $p \in I^+(S)$, contradicting that it is on the boundary. It follows that $\gamma(t) \in \partial I^+(S)$. Thus we have shown that $\gamma|_{[0,a)} \subset \partial I^+(S)$. Suppose for the moment $\gamma|_{[0,a)}$ is piecewise smooth. Since $\partial I^+(S)$ is achronal, no two points of γ can be joined by a timelike curve. It then follows from Proposition 2.1 that γ is a null geodesic. But using the fact that C^0 causal curves can be approximated by piecewise smooth causal curves, one can show in the general case that $\gamma|_{[0,a)}$ is a null geodesic. (Exercise: Show this.)

Finally, we consider the two cases $a < \infty$ and $a = \infty$. If $a < \infty$, then by the uniform convergence, $\gamma(a) = \lim_{m \rightarrow \infty} \gamma_m(a_m) = \lim_{m \rightarrow \infty} q_m \in S$, since S is closed. Thus, we have a null geodesic from p contained in $\partial I^+(S)$ that ends on S . If $a = \infty$ then we have a null geodesic from p in $\partial I^+(S)$ that is past inextendible in M . \square

Achronal boundaries have been recently employed in a fundamental way to study Lorentzian splitting problems, cf. [9].

4 Causality conditions

A number of results in Lorentzian geometry and general relativity require some sort of causality condition. It is perhaps natural on physical grounds to rule out the occurrence of closed timelike curves. Physically, the existence of such a curve signifies the existence of an observer who is able to travel into his/her own past, which leads to variety of paradoxical situations. A spacetime M satisfies the *chronology condition* provided there are no closed timelike curves in M . Compact spacetimes have limited interest in general relativity since they all violate the chronology condition.

Proposition 4.1. *Every compact spacetime contains a closed timelike curve.*

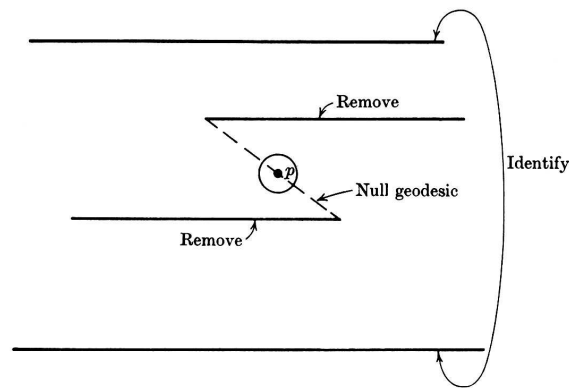
Proof. The sets $\{I^+(p); p \in M\}$ form an open cover of M from which we can abstract a finite subcover: $I^+(p_1), I^+(p_2), \dots, I^+(p_k)$. We may assume that this is the minimal number of such sets covering M . Since these sets cover M , $p_1 \in I^+(p_i)$ for some i . It follows that $I^+(p_1) \subset I^+(p_i)$. Hence, if $i \neq 1$, we could reduce the number of sets

in the cover. Thus, $p_1 \in I^+(p_1)$ which implies that there is a closed timelike curve through p_1 . \square

A somewhat stronger condition than the chronology condition is the *causality condition*. A spacetime M satisfies the causality condition provided there are no closed (nontrivial) causal curves in M .

Exercise: Construct a spacetime that satisfies the chronology condition but not the causality condition.

A spacetime that satisfies the causality condition can nonetheless be on the verge of failing it, in the sense that there exist causal curves that are “almost closed”, as illustrated by the following figure.



Strong causality is a condition that rules out almost closed causal curves. An open set U in spacetime M is said to be *causally convex* provided no causal curve in M meets U in a disconnected set. Given $p \in M$, strong causality is said to hold at p provided p has arbitrarily small causally convex neighborhoods, i.e., for each neighborhood V of p there exists a causally convex neighborhood U of p such that $U \subset V$. Note that strong causality fails at the point p in the figure above. In fact strong causality fails at all points along the dashed null geodesic. It can be shown that the set of points at which strong causality holds is open.

M is said to be strongly causal if strong causality holds at all of its points. This is the “standard” causality condition in spacetime geometry, and, although there are even stronger causality conditions, it is sufficient for most applications. There is an interesting connection between strong causality and the so-called *Alexandrov topology*. The sets of the form $I^+(p) \cap I^-(q)$ form the base for a topology on M , which is the Alexandrov topology. This topology is in general more coarse than the manifold topology of M . However it can be shown that the Alexandrov topology agrees with the manifold topology iff M is strongly causal.

The following lemma is often useful.

Lemma 4.2. *Suppose strong causality holds at each point of a compact set K in a spacetime M . If $\gamma : [0, b) \rightarrow M$ is a future inextendible causal curve that starts in K*

then eventually it leaves K and does not return, i.e., there exists $t_0 \in [0, b)$ such that $\gamma(t) \notin K$ for all $t \in [t_0, b)$.

Proof. Exercise.

In referring to the property described by this lemma, we say that a future inextendible causal curve cannot be “imprisoned” or “partially imprisoned” in a compact set on which strong causality holds.

We now come to a fundamental condition in spacetime geometry, that of *global hyperbolicity*. Mathematically, global hyperbolicity is a basic ‘niceness’ condition that often plays a role analogous to geodesic completeness in Riemannian geometry. Physically, global hyperbolicity is connected to the notion of (strong) cosmic censorship introduced by Roger Penrose. This is the conjecture that, generically, spacetime solutions to the Einstein equations do not admit *naked singularities* (singularities visible to some observer).

Definition 4.1. *A spacetime M is said to be globally hyperbolic provided*

- (1) M is strongly causal.
- (2) (*Internal Compactness*) The sets $J^+(p) \cap J^-(q)$ are compact for all $p, q \in M$.

Condition (2) says roughly that M has no holes or gaps. For example Minkowski space is globally hyperbolic but the spacetime obtained by removing one point from it is not. We note that it can be shown that the causality condition, together with internal compactness, implies strong causality, so that causality could replace strong causality in the definition of global hyperbolicity; cf. [3]. However, causality alone is not sufficient to guarantee the conclusion of Lemma 4.2.

We consider a few basic consequences of global hyperbolicity.

Proposition 4.3. *Let M be a globally hyperbolic spacetime. Then,*

- (1) The sets $J^\pm(A)$ are closed, for all compact $A \subset M$.
- (2) The sets $J^+(A) \cap J^-(B)$ are compact, for all compact $A, B \subset M$.

Proof. We prove $\overline{J^\pm(p)}$ are closed for all $p \in M$, and leave the rest as an exercise. Suppose $q \in \overline{J^+(p)} \setminus J^+(p)$ for some $p \in M$. Choose $r \in I^+(q)$, and $\{q_n\} \subset J^+(p)$, with $q_n \rightarrow q$. Since $I^-(r)$ is an open neighborhood of q , $\{q_n\} \subset J^-(r)$ for n large. It follows that $q \in \overline{J^+(p) \cap J^-(r)} = J^+(p) \cap J^-(r)$, since $J^+(p) \cap J^-(r)$ is compact and hence closed. But this contradicts $q \notin J^+(p)$. Thus, $J^+(p)$ is closed, and similarly so is $J^-(p)$. \square

Analogously to the case of Riemannian geometry, one can learn much about the global structure of spacetime by studying its causal geodesics. Locally, causal geodesics maximize Lorentzian arc length (cf., Proposition 1.3). Given $p, q \in M$, with $p < q$, we wish to consider conditions under which there exists a maximal future

directed causal geodesic γ from p to q , where by maximal we mean that for any future directed causal curve σ from p to q , $L(\gamma) \geq L(\sigma)$.

Maximality can be conveniently expressed in terms of the Lorentzian *distance function*, $d : M \times M \rightarrow [0, \infty]$. For $p < q$, let $\Omega_{p,q}$ denote the collection of future directed causal curves from p to q . Then, for any $p, q \in M$, define

$$d(p, q) = \begin{cases} \sup\{L(\sigma) : \sigma \in \Omega_{p,q}\}, & \text{if } p < q \\ 0, & \text{if } p \not< q \end{cases}$$

While the Lorentzian distance function is not a distance function in the usual sense of metric spaces, and may not even be finite valued, it does have a few nice properties. For one, it obeys a *reverse triangle inequality*,

$$\text{if } p < r < q \quad \text{then} \quad d(p, q) \geq d(p, r) + d(r, q).$$

Exercise: Prove this.

We have the following basic fact.

Proposition 4.4. *The Lorentzian distance function is lower semi-continuous.*

Proof. Fix $p, q \in M$. Given $\epsilon > 0$ we need to find neighborhoods U and V of p and q , respectively, such that for all $x \in U$ and all $y \in V$, $d(x, y) > d(p, q) - \epsilon$.

If $d(p, q) = 0$ there is nothing to prove. Thus, we assume $p < q$ and $0 < d(p, q) < \infty$. We leave the case $d(p, q) = \infty$ as an exercise. Let σ be a future directed timelike curve from p to q such that $L(\sigma) = d(p, q) - \epsilon/3$. Let U and V be convex neighborhoods of p and q , respectively. Choose p' on σ close to p and q' on σ close to q . Then $U' = I^-(p', U)$ and $V' = I^+(q', V)$ are neighborhoods of p and q , respectively. Moreover, by choosing p' sufficiently close to p and q' sufficiently close to q , one verifies that for all $x \in U'$ and $y \in V'$, there exists a future directed timelike curve α from x to y , containing the portion of σ from p' to q' , having length $L(\alpha) > d(p, q) - \epsilon/2$. \square

Though the Lorentzian distance function is not continuous in general, it is continuous (and finite valued) for globally hyperbolic spacetimes; cf., [15, p. 412].

Given $p < q$, note that a causal geodesic segment γ having length $L(\gamma) = d(p, q)$ is maximal. Global hyperbolicity is the standard condition to ensure the existence of maximal causal geodesic segments.

Proposition 4.5. *Let M be a globally hyperbolic spacetime. Given $p, q \in M$, with $q \in I^+(p)$, there exists a maximal future directed timelike geodesic γ from p to q ($L(\gamma) = d(p, q)$).*

Proof. The proof involves a standard limit curve argument, together with the fact that the Lorentzian arc length functional is upper semi-continuous; see [17, p. 54].

As usual, let h be a background complete Riemannian metric on M . For each n , let $\gamma_n : [0, a_n] \rightarrow M$ be a future directed timelike curve from $p = \gamma_n(0)$ to $q = \gamma_n(a_n)$, parameterized with respect to h -arc length, such that $L(\gamma_n) \rightarrow d(p, q)$. Extend each γ_n

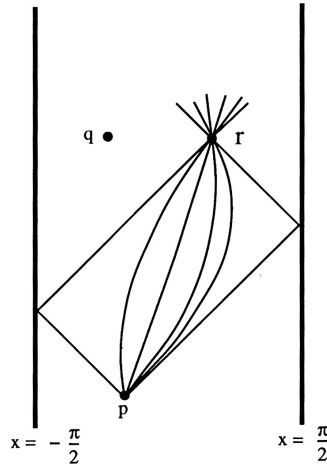
to a future inextendible causal curve $\tilde{\gamma}_n : [0, \infty) \rightarrow M$, parameterized with respect to h -arc length. By the limit curve lemma, there exists a subsequence $\tilde{\gamma}_m : [0, \infty) \rightarrow M$ that converges to a future inextendible C^0 causal curve $\gamma : [0, \infty) \rightarrow M$ such that $\gamma(0) = p$. By taking a further subsequence if necessary we can assume $a_m \uparrow a$. Since each γ_m is contained in the compact set $J^+(p) \cap J^-(q)$, it follows that $\gamma|_{[0,a]} \subset \overline{J^+(p) \cap J^-(q)} = J^+(p) \cap J^-(q)$. Since M is strongly causal, it must be that $a < \infty$, otherwise, γ would be imprisoned in $J^+(p) \cap J^-(q)$, contradicting Lemma 4.2. Then, $\gamma(a) = \lim_{m \rightarrow \infty} \gamma_m(a_m) = q$.

Let $\bar{\gamma} = \gamma|_{[0,a]}$. $\bar{\gamma}$ is a future directed C^0 causal curve from p to q . Moreover, by the upper semi-continuity of L ,

$$L(\bar{\gamma}) \geq \limsup_{m \rightarrow \infty} L(\gamma_m) = d(p, q),$$

and so $L(\bar{\gamma}) = d(p, q)$. Hence, $\bar{\gamma}$ has maximal length among all future directed causal curves from p to q . This forces each sub-segment of $\bar{\gamma}$ to have maximal length. Using Proposition 1.3 (part (2) of which remains valid for C^0 causal curves) and Proposition 2.1, one can then argue that each sufficiently small segment of $\bar{\gamma}$ is a causal geodesic. (Exercise: Argue this.) \square

Remarks: There are simple examples showing that if either of the conditions (1) or (2) fail to hold in the definition of global hyperbolicity then maximal segments may fail to exist. Moreover, contrary to the situation in Riemannian geometry, geodesic completeness does not guarantee the existence of maximal segments, as is well illustrated by anti-de Sitter space which is geodesically complete. The figure below depicts 2-dimensional anti-de Sitter space. It can be represented as the strip $M = \{(t, x) : -\pi/2 < x < \pi/2\}$, equipped with the metric $ds^2 = \sec^2 x(-dt^2 + dx^2)$. Because the anti-de Sitter metric is conformal to the Minkowski metric on the strip, pasts and futures of both space times are the same. It can be shown that all future directed timelike geodesics emanating from p refocus at r . The points p and q are timelike related, but there is no timelike geodesic segment from p to q .



Global hyperbolicity is closely related to the existence of certain ‘ideal initial value hypersurfaces’, called *Cauchy surfaces*. There are slight variations in the literature in the definition of a Cauchy surface. Here we adopt the following definition.

Definition 4.2. *A Cauchy surface for a spacetime M is an achronal subset S of M which is met by every inextendible causal curve in M .*

From the definition one can easily show that if S is a Cauchy surface for M then $S = \partial I^+(S) = \partial I^-(S)$ (exercise!). It follows from Proposition 3.1 that a Cauchy surface S is a closed achronal C^0 hypersurface in M .

Theorem 4.6 (Geroch, [11]). *If a spacetime M is globally hyperbolic then it has a Cauchy surface S , and conversely.*

Proof. We make a couple brief comments about the proof. The converse will be discussed in the next section. To prove that a globally hyperbolic spacetime M admits a Cauchy surface, one introduces a measure μ on M such that $\mu(M) = 1$. Consider the function $f : M \rightarrow \mathbb{R}$ defined by

$$f(p) = \frac{\mu(J^-(p))}{\mu(J^+(p))}.$$

Internal compactness is used to show that f is continuous, and strong causality is used to show that f is strictly increasing along future directed causal curves. Moreover, if $\gamma : (a, b) \rightarrow M$ is a future directed inextendible causal curve in M , one shows $f(\gamma(t)) \rightarrow 0$ as $t \rightarrow a^+$, and $f(\gamma(t)) \rightarrow \infty$ as $t \rightarrow b^-$. It follows that $S = \{p \in M : f(p) = 1\}$ is a Cauchy surface for M . \square

Remark: The function f constructed in the proof is what is referred to as a *time function*, namely, a continuous function that is strictly increasing along future directed causal curves. See e.g, [2, 7] for recent developments concerning the construction of *smooth time functions*, i.e., smooth functions with timelike gradient (which hence are necessarily time functions) and smooth spacelike Cauchy surfaces.

Proposition 4.7. *Let M be globally hyperbolic.*

- (1) *If S is a Cauchy surface for M then M is homeomorphic to $\mathbb{R} \times S$.*
- (2) *Any two Cauchy surfaces in M are homeomorphic.*

Proof. To prove (1), one introduces a future directed timelike vector field X on M . X can be scaled so that the time parameter t of each integral curve of X extends from $-\infty$ to ∞ , with $t = 0$ at points of S . Each $p \in M$ is on an integral curve of X that meets S in exactly one point q . This sets up a correspondence $p \leftrightarrow (t, q)$, which gives the desired homeomorphism. A similar technique may be used to prove (2) \square

In view of Proposition 4.7, any nontrivial topology in a globally hyperbolic spacetime must reside in its Cauchy surfaces.

The following fact is often useful.

Proposition 4.8. *If S is a compact achronal C^0 hypersurface in a globally hyperbolic spacetime M then S must be a Cauchy surface for M .*

The proof will be discussed in the next section.

5 Domains of dependence

Definition 5.1. *Let S be an achronal set in a spacetime M . We define the future and past domains of dependence of S , $D^+(S)$ and $D^-(S)$, respectively, as follows,*

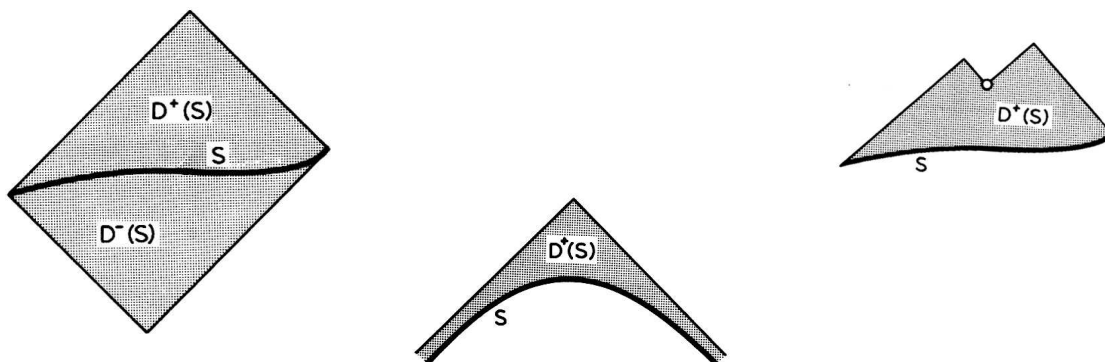
$$D^+(S) = \{p \in M : \text{every past inextendible causal curve from } p \text{ meets } S\},$$

$$D^-(S) = \{p \in M : \text{every future inextendible causal curve from } p \text{ meets } S\}.$$

The (total) domain of dependence of S is the union, $D(S) = D^+(S) \cup D^-(S)$.

In physical terms, since information travels along causal curves, a point in $D^+(S)$ only receives information from S . Thus if physical laws are suitably causal, initial data on S should determine the physics on $D^+(S)$ (in fact on all of $D(S)$).

Below we show a few examples of future and past domains of dependence.



The figure in the right shows the effect of removing a point from M . The figure in the center shows the future domain of dependence of the spacelike hyperboloid $t^2 - x^2 = 1$, $t < 0$, in the Minkowski plane.

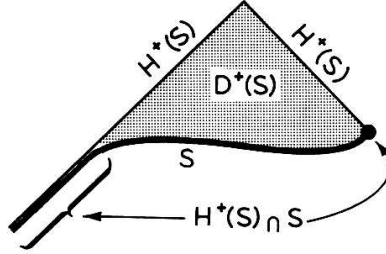
If S is achronal, the future Cauchy horizon $H^+(S)$ of S is the future boundary of $D^+(S)$. This is made precise in the following definition.

Definition 5.2. *Let $S \subset M$ be achronal. The future Cauchy horizon $H^+(S)$ of S is defined as follows*

$$H^+(S) = \{p \in \overline{D^+(S)} : I^+(p) \cap D^+(S) = \emptyset\}$$

$$= \overline{D^+(S)} \setminus I^-(D^+(S)).$$

The past Cauchy horizon $H^-(S)$ is defined time-dually. The (total) Cauchy horizon of S is defined as the union, $H(S) = H^+(S) \cup H^-(S)$. (See the figure below.)



We record some basic facts about domains of dependence and Cauchy horizons.

Proposition 5.1. *Let S be an achronal subset of M . Then the following hold.*

- (1) $S \subset D^+(S)$.
- (2) If $p \in D^+(S)$ then $I^-(p) \cap I^+(S) \subset D^+(S)$.
- (3) $\partial D^+(S) = H^+(S) \cup S$.
- (4) $H^+(S)$ is achronal.
- (5) $\text{edge } H^+(S) \subset \text{edge } S$, with equality holding if S is closed.

(4): The achronality of $H^+(S)$ follows almost immediately from the definition: Suppose $p, q \in H^+(S)$ with $p \ll q$. Since $q \in \overline{D^+(S)}$, and $I^+(p)$ is a neighborhood of q , $I^+(p)$ meets $D^+(S)$, contradicting the definition of $H^+(S)$. We leave the proofs of the other parts as an exercise.

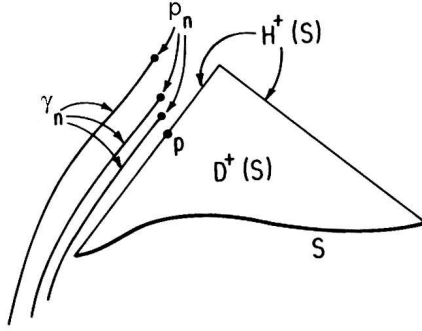
Cauchy horizons have structural properties similar to achronal boundaries, as indicated in the next two results. From Proposition 3.3 we obtain the following.

Proposition 5.2. *Let $S \subset M$ be achronal. Then $H^+(S) \setminus \text{edge } H^+(S)$, if nonempty, is an achronal C^0 hypersurface in M .*

In a similar vein to Proposition 3.4, we have the following.

Proposition 5.3. *Let S be an achronal subset of M . Then $H^+(S)$ is ruled by null geodesics, i.e., every point of $H^+(S) \setminus \text{edge } S$ is the future endpoint of a null geodesic in $H^+(S)$ which is either past inextendible in M or else has a past end point on edge S .*

Comments on the proof. The proof uses a limit curve argument. Consider the case $p \in H^+(S) \setminus S$. Since $I^+(p) \cap D^+(S) = \emptyset$, we can find a sequence of points $p_n \notin D^+(S)$, such that $p_n \rightarrow p$. For each n , there exists a past inextendible causal curve γ_n that does not meet S . By the limit curve lemma there exists a subsequence γ_m that converges to a past inextendible C^0 causal curve γ starting at p . Near p this defines the desired null geodesic (see the figure below).



The case $p \in S \setminus \text{edge } S$ is handled somewhat differently; for details see [18, p. 203]. \square

The next proposition follows straight-forwardly from definitions, together with the fact that, for S achronal, $\partial D(S) = H(S)$ (exercise).

Proposition 5.4. *Let S be an achronal subset of a spacetime M . Then, S is a Cauchy surface for M if and only if $D(S) = M$ if and only if $H(S) = \emptyset$.*

The following basic result ties domains of dependence to global hyperbolicity.

Proposition 5.5. *Let $S \subset M$ be achronal.*

- (1) *Strong causality holds on $\text{int } D(S)$.*
- (2) *Internal compactness holds on $\text{int } D(S)$, i.e., for all $p, q \in \text{int } D(S)$, $J^+(p) \cap J^-(q)$ is compact.*

Comments on the proof. With regard to (1), first observe that the chronology condition holds on $D(S)$. For instance, suppose there exists a timelike curve γ passing through $p \in D^+(S)$, and take γ to be past directed. By repeating loops we obtain a past inextendible timelike curve $\tilde{\gamma}$, which hence must meet S . In fact, it will meet S infinitely often, thereby violating the achronality of S . A similar argument shows that the causality condition holds on $\text{int } D(S)$. Suppose for example that γ is a past directed closed causal curve through $p \in \text{int } D^+(S)$. By repeating loops we obtain a past inextendible causal curve $\tilde{\gamma}$ starting at p . Thus $\tilde{\gamma}$ meets S , and since $p \in \text{int } D^+(S)$, will enter $I^-(S)$ (see Lemma 5.8 below). This again leads to an achronality violation. By more refined arguments, using the limit curve lemma, one can show that strong causality holds on $\text{int } D(S)$. With regard to (2), suppose there exist $p, q \in \text{int } D(S)$, such that $J^+(p) \cap J^-(q)$ is noncompact. We want to show that every sequence q_n in $J^+(p) \cap J^-(q)$ has a convergent subsequence. Without loss of generality we may assume $\{q_n\} \subset D^-(S)$. For each n , let γ_n be a future directed causal curve from p to q passing through q_n . As usual, extend each γ_n to a future inextendible causal curve $\tilde{\gamma}_n$. By the limit curve lemma, there exists a subsequence $\tilde{\gamma}_m$ that converges to a future inextendible C^0 causal curve γ starting at p . One can then show that either the sequence of points q_m converges or γ does not enter $I^+(S)$. \square

We can now address the converse part of Theorem 4.6.

Corollary 5.6. *If S is a Cauchy surface for M then M is globally hyperbolic.*

Proof. This follows immediately from Propositions 5.4 and 5.5: S Cauchy $\implies D(S) = M \implies \text{int } D(S) = M \implies M$ is globally hyperbolic. \square

We now give a proof of Proposition 4.8 from the previous section.

Proof of Proposition 4.8. It suffices to show that $H(S) = H^+(S) \cup H^-(S) = \emptyset$. Suppose there exists $p \in H^+(S)$. Since S is edgeless, it follows from Proposition 5.3 that p is the future endpoint of a past inextendible null geodesic $\gamma \subset H^+(S)$. Then since $\gamma \subset D^+(S) \cap J^-(p)$ (exercise: show this), we have that γ is contained in the set $J^+(S) \cap J^-(p)$, which is compact by Proposition 4.3. By Lemma 4.2 strong causality must be violated at some point of $J^+(S) \cap J^-(p)$. Thus $H^+(S) = \emptyset$, and time-dually, $H^-(S) = \emptyset$.

We conclude this section by stating several lemmas that are useful in proving some of the results described here, as well as other results concerning domains of dependence.

Lemma 5.7 ([15], p. 416). *Let γ be a past inextendible causal curve starting at p that does not meet a closed set C . If $p_0 \in I^+(p, M \setminus C)$, there exists a past inextendible timelike curve starting at p_0 that does not meet C .*

Proof. Exercise.

Lemma 5.8. *Let S be achronal. If $p \in \text{int } D^+(S)$ then every past inextendible causal curve from p enters $I^-(S)$.*

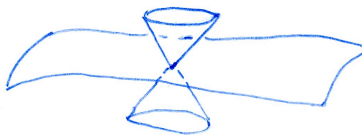
Proof. This follows from the *proof* of the preceding lemma.

Lemma 5.9. *Let S be achronal. Then $p \in \overline{D^+(S)}$ iff every past inextendible timelike curve meets S .*

Proof. Exercise.

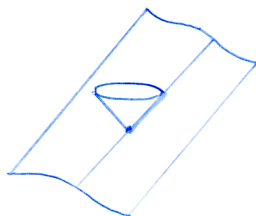
6 The geometry of null hypersurfaces

In addition to curves, one can discuss the causality of certain higher dimensional submanifolds. For example, a *spacelike* hypersurface is a hypersurface all of whose tangent vectors are spacelike, or, equivalently, whose normal vectors are timelike:



More precisely, a hypersurface is spacelike if the induced metric is positive definite (i.e. Riemannian). In GR, a spacelike hypersurface represents space at a given instant of time.

A null hypersurface is a hypersurface such that the null cone is tangent to at each of its points:



Null hypersurfaces play an important role in GR as they represent horizons of various sorts. For example the event horizons in the Schwarzschild and Kerr spacetimes are null hypersurfaces. Null hypersurfaces have an interesting geometry which we would like to discuss in this section.

In more precise terms a null hypersurface in a spacetime (M, g) is a smooth co-dimension one submanifold S of M , such that at each $p \in M$, $g : T_p S \times T_p S \rightarrow \mathbb{R}$ is degenerate. This means that there exists a nonzero vector $K_p \in T_p S$ (the direction of degeneracy) such that

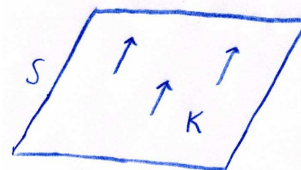
$$\langle K_p, X \rangle = 0 \quad \text{for all } X \in T_p S$$

where we have introduced the shorthand metric notation: $\langle U, V \rangle = g(U, V)$. In particular,

- (1) K_p is a null vector, $\langle K_p, K_p \rangle = 0$, which we can choose to be future pointing, and
- (2) $[K_p]^\perp = T_p S$.
- (3) Moreover, every (nonzero) vector $X \in T_p S$ that is not a multiple of K_p is spacelike.

Thus, every null hypersurface S gives rise to a future directed null vector field K ,

$$p \in S \xrightarrow{K} K_p \in T_p S,$$



which will be smooth, $K \in \mathfrak{X}(S)$, provided it is normalized in a suitably uniform way. Furthermore, the null vector field K is unique up to a positive pointwise scale factor.

As simple examples, in Minkowski space \mathbb{M}^{n+1} , the past and future cones, $\partial I^-(p)$ and $\partial I^+(p)$, respectively, are smooth null hypersurfaces away from the vertex p . Each

nonzero null vector $X \in T_p\mathbb{M}^{n+1}$ determines a null hyperplane $\Pi = \{q \in \mathbb{M}^{n+1} : \langle \overline{pq}, X \rangle = 0\}$.

The following fact is fundamental.

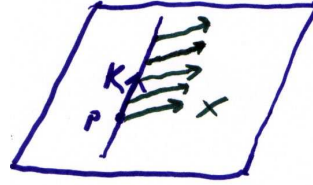
Proposition 6.1. *Let S be a smooth null hypersurface and let $K \in \mathfrak{X}(S)$ be a smooth future directed null vector field on S . Then the integral curves of K are null geodesics (when suitably parameterized),*

Remark: The integral curves of K are called the null generators of S . Apart from parameterizations, the null generators are intrinsic to the null hypersurface.

Proof. It suffices to show that $\nabla_K K = \lambda K$, for then the integral curves are in general *pre-geodesics* (i.e., are geodesics after a suitable reparameterization). To show this it suffice to show that at each $p \in S$, $\nabla_K K \perp T_p S$, i.e., $\langle \nabla_K K, X \rangle = 0$ for all $X \in T_p S$.

Extend $X \in T_p S$ by making it invariant under the flow generated by K ,

$$[K, X] = \nabla_K X - \nabla_X K = 0$$



X remains tangent to S , so along the flow line through p , $\langle K, X \rangle = 0$. Differentiating we obtain,

$$0 = K\langle K, X \rangle = \langle \nabla_K K, X \rangle + \langle K, \nabla_K X \rangle,$$

and hence,

$$\langle \nabla_K K, X \rangle = -\langle K, \nabla_K X \rangle = -\langle K, \nabla_X K \rangle = -\frac{1}{2}X\langle K, K \rangle = 0.$$

□

To study the ‘shape’ of the null hypersurface S we study how the null vector field K varies along S . Since K is actually orthogonal to S , this is somewhat analogous to how we study the shape of a hypersurface in a Riemannian manifold, or spacelike hypersurface in a Lorentzian manifold, by introducing the shape operator (or Weingarten map) and associated second fundamental form. We proceed to introduce null analogues of these objects. For technical reasons one works “mod K ”, as described below.

We introduce the following equivalence relation on tangent vectors: For $X, X' \in T_p S$,

$$X' = X \text{ mod } K \quad \text{if and only if} \quad X' - X = \lambda K \text{ for some } \lambda \in \mathbb{R}.$$

Let \overline{X} denote the equivalence class of X . Let $T_p S/K = \{\overline{X} : X \in T_p S\}$, and $TS/K = \cup_{p \in S} T_p S/K$. TS/K , the mod K tangent bundle of S , is a smooth rank $n - 1$ vector bundle over S . This vector bundle does not depend on the particular choice of null vector field K .

There is a natural positive definite metric h on TS/K induced from $\langle \cdot, \cdot \rangle$: For each $p \in S$, define $h : T_p S/K \times T_p S/K \rightarrow \mathbb{R}$ by $h(\bar{X}, \bar{Y}) = \langle X, Y \rangle$. A simple computation shows that h is well-defined: If $X' = X \bmod K$, $Y' = Y \bmod K$ then

$$\begin{aligned} \langle X', Y' \rangle &= \langle X + \alpha K, Y + \beta K \rangle \\ &= \langle X, Y \rangle + \beta \langle X, K \rangle + \alpha \langle K, Y \rangle + \alpha \beta \langle K, K \rangle \\ &= \langle X, Y \rangle. \end{aligned}$$

The *null Weingarten map* $b = b_K$ of S with respect to K is, for each point $p \in S$, a linear map $b : T_p S/K \rightarrow T_p S/K$ defined by $b(\bar{X}) = \overline{\nabla_X K}$.

Exercise: Show that b is well-defined. Show also that that if $\tilde{K} = fK$, $f \in C^\infty(S)$, is any other future directed null vector field on S , then $b_{\tilde{K}} = fb_K$. It follows that the Weingarten map $b = b_K$ at a point p is uniquely determined by the value of K at p .

Proposition 6.2. b is self adjoint with respect to h , i.e., $h(b(\bar{X}), \bar{Y}) = h(\bar{X}, b(\bar{Y}))$, for all $\bar{X}, \bar{Y} \in T_p S/K$.

Proof. Extend $X, Y \in T_p S$ to vector fields tangent to S near p . Using $X \langle K, Y \rangle = 0$ and $Y \langle K, X \rangle = 0$, we obtain,

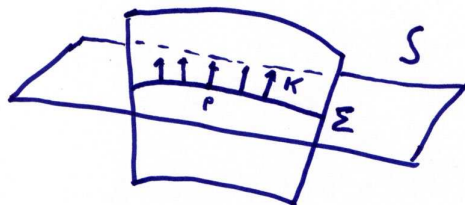
$$\begin{aligned} h(b(\bar{X}), \bar{Y}) &= \langle \nabla_X K, Y \rangle = -\langle K, \nabla_X Y \rangle = -\langle K, \nabla_Y X \rangle + \langle K, [X, Y] \rangle \\ &= \langle \nabla_Y K, X \rangle = h(\bar{X}, b(\bar{Y})). \end{aligned}$$

The *null second fundamental form* $B = B_K$ of S with respect to K is the bilinear form associated to b via h : For each $p \in S$, $B : T_p S/K \times T_p S/K \rightarrow \mathbb{R}$ is defined by,

$$B(\bar{X}, \bar{Y}) = h(b(\bar{X}), \bar{Y}) = \langle \nabla_X K, Y \rangle.$$

Since b is self-adjoint, B is symmetric. We say that S is *totally geodesic* iff $B \equiv 0$. This has the usual geometric meaning: If S is totally geodesic then any geodesic in M starting tangent to S stays in S . This follows from the fact that, when S is totally geodesic, the restriction to S of the Levi-Civita connection of M defines a linear connection on S . Null hyperplanes in Minkowski space are totally geodesic, as is the event horizon in Schwarzschild spacetime.

The *null mean curvature* (or *null expansion scalar*) of S with respect to K is the smooth scalar field θ on S defined by, $\theta = \text{tr } b$. θ has a natural geometric interpretation. Let Σ be the intersection of S with a hypersurface in M which is transverse to K near $p \in S$; Σ will be a co-dimension two spacelike submanifold of M , along which K is orthogonal.



Let $\{e_1, e_2, \dots, e_{n-1}\}$ be an orthonormal basis for $T_p\Sigma$ in the induced metric. Then $\{\bar{e}_1, \bar{e}_2, \dots, \bar{e}_{n-1}\}$ is an orthonormal basis for T_pS/K . Hence at p ,

$$\begin{aligned}\theta &= \text{tr } b = \sum_{i=1}^{n-1} h(b(\bar{e}_i), \bar{e}_i) = \sum_{i=1}^{n-1} \langle \nabla_{e_i} K, e_i \rangle. \\ &= \text{div}_\Sigma K.\end{aligned}\tag{6.15}$$

where $\text{div}_\Sigma K$ is the divergence of K along Σ . Thus, θ measures the overall expansion of the null generators of S towards the future.

It follows from the exercise on the preceding page that if $\tilde{K} = fK$ then $\tilde{\theta} = f\theta$. Thus the null mean curvature inequalities $\theta > 0$, $\theta < 0$, etc., are invariant under positive rescaling of K . In Minkowski space, a future null cone $S = \partial I^+(p) \setminus \{p\}$ (resp., past null cone $S = \partial I^-(p) \setminus \{p\}$) has positive null mean curvature, $\theta > 0$ (resp., negative null mean curvature, $\theta < 0$).

We now study how the null Weingarten map propagates along the null geodesic generators of S . Let $\eta : I \rightarrow M$, $s \rightarrow \eta(s)$, be a future directed affinely parameterized null geodesic generator of S . For each $s \in I$, let

$$b(s) = b_{\eta'(s)} : T_{\eta(s)}S/\eta'(s) \rightarrow T_{\eta(s)}S/\eta'(s)\tag{6.16}$$

be the Weingarten map based at $\eta(s)$ with respect to the null vector $K = \eta'(s)$. We show that the one parameter family of Weingarten maps $s \rightarrow b(s)$, obeys a certain Riccati equation.

We first need to make a few definitions. Let $s \rightarrow \mathcal{Y}(s)$ be a TS/K vector field along η , i.e., for each $s \in I$, $\mathcal{Y}(s) \in T_{\eta(s)}S/K$. We say that $s \rightarrow \mathcal{Y}(s)$ is smooth if, at least locally, there is a smooth (in the usual sense) vector field $s \rightarrow Y(s)$ along η , tangent to S , such that $\mathcal{Y}(s) = \overline{Y(s)}$. Then define the covariant derivative of $s \rightarrow \mathcal{Y}(s)$ along η by, $\mathcal{Y}'(s) = \overline{Y'(s)}$, where Y' is the usual covariant differentiation.

Exercise: Show that \mathcal{Y}' is independent of the choice of Y .

Then the covariant derivative of b along η is defined by requiring a natural product rule to hold. If $s \rightarrow X(s)$ is a vector field along η tangent to S , b' is defined by,

$$b'(\overline{X}) = b(\overline{X})' - b(\overline{X}').\tag{6.17}$$

Proposition 6.3. *The one parameter family of Weingarten maps $s \rightarrow b(s)$, obeys the following Riccati equation,*

$$b' + b^2 + R = 0,\tag{6.18}$$

where $R : T_{\eta(s)}S/\eta'(s) \rightarrow T_{\eta(s)}S/\eta'(s)$ is the curvature endomorphism defined by $R(\overline{X}) = \overline{R(X, \eta'(s))\eta'(s)}$.

Proof. Fix a point $p = \eta(s_0)$, $s_0 \in (a, b)$, on η . On a neighborhood U of p in S we can scale the null vector field K so that K is a geodesic vector field, $\nabla_K K = 0$, and so that K , restricted to η , is the velocity vector field to η , i.e., for each s near s_0 , $K_{\eta(s)} = \eta'(s)$. Let $X \in T_p M$. Shrinking U if necessary, we can extend X to a smooth vector field on U so that $[X, K] = \nabla_X K - \nabla_K X = 0$. Then,

$$R(X, K)K = \nabla_X \nabla_K K - \nabla_K \nabla_X K - \nabla_{[X, K]} K = -\nabla_K \nabla_K X.$$

Hence along η we have, $X'' = -R(X, \eta')\eta'$ (which implies that X , restricted to η , is a *Jacobi field* along η). Thus, from Equation 6.17, at the point p we have,

$$\begin{aligned} b'(\bar{X}) &= \overline{\nabla_X K}' - b(\overline{\nabla_K X}) = \overline{\nabla_K X}' - b(\overline{\nabla_X K}) \\ &= \overline{X''} - b(b(\bar{X})) = -\overline{R(X, \eta')\eta'} - b^2(\bar{X}) \\ &= -R(\bar{X}) - b^2(\bar{X}), \end{aligned}$$

which establishes Equation 6.18. \square

By taking the trace of (6.18) we obtain the following formula for the derivative of the null mean curvature $\theta = \theta(s)$ along η ,

$$\theta' = -\text{Ric}(\eta', \eta') - \sigma^2 - \frac{1}{n-1}\theta^2, \quad (6.19)$$

where $\sigma := (\text{tr } \hat{b}^2)^{1/2}$ is the *shear scalar*, $\hat{b} := b - \frac{1}{n-1}\theta \cdot \text{id}$ is the trace free part of the Weingarten map, and $\text{Ric}(\eta', \eta') = R_{ij}(\eta^i)'(\eta^j)'$ is the Ricci tensor contracted on the tangent vector η' . Equation 6.19 is known in relativity as the Raychaudhuri equation (for an irrotational null geodesic congruence). This equation shows how the Ricci curvature of spacetime influences the null mean curvature of a null hypersurface.

The following proposition is a standard application of the Raychaudhuri equation.

Proposition 6.4. *Let M be a spacetime which obeys the null energy condition (NEC), $\text{Ric}(X, X) \geq 0$ for all null vectors X , and let S be a smooth null hypersurface in M . If the null generators of S are future geodesically complete then S has nonnegative null mean curvature, $\theta \geq 0$.*

Proof. Suppose $\theta < 0$ at $p \in S$. Let $s \rightarrow \eta(s)$ be the null generator of S passing through $p = \eta(0)$, affinely parametrized. Let $b(s) = b_{\eta'(s)}$, and take $\theta = \text{tr } b$. By the invariance of sign under scaling, one has $\theta(0) < 0$. Raychaudhuri's equation and the NEC imply that $\theta = \theta(s)$ obeys the inequality,

$$\frac{d\theta}{ds} \leq -\frac{1}{n-1}\theta^2, \quad (6.20)$$

and hence $\theta < 0$ for all $s > 0$. Dividing through by θ^2 then gives,

$$\frac{d}{ds} \left(\frac{1}{\theta} \right) \geq \frac{1}{n-1}, \quad (6.21)$$

which implies $1/\theta \rightarrow 0$, i.e., $\theta \rightarrow -\infty$ in finite affine parameter time, contradicting the smoothness of θ . \square

Exercise. Let Σ be a local cross section of the null hypersurface S , as depicted on p. 26, with volume form ω . If Σ is moved under flow generated by K , show that $L_K\omega = \theta\omega$, where $L = \text{Lie derivative}$.

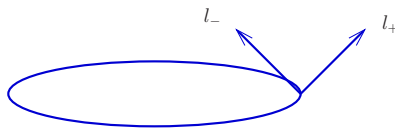
Thus, Proposition 6.4 implies, under the given assumptions, that cross sections of S are nondecreasing in area as one moves towards the future. Proposition 6.4 is the simplest form of Hawking's black hole area theorem [12]. For a study of the area theorem, with a focus on issues of regularity, see [6].

7 Trapped surfaces and the Penrose Singularity Theorem

In this section we introduce the important notion of a trapped surface and present the classical Penrose singularity theorem.

Let (M^{n+1}, g) be an $(n + 1)$ -dimensional spacetime, with $n \geq 3$. Let Σ^{n-1} be a closed (i.e., compact without boundary) co-dimension two spacelike submanifold of M . Each normal space of Σ , $[T_p\Sigma]^\perp$, $p \in \Sigma$, is timelike and 2-dimensional, and hence admits two future directed null directions orthogonal to Σ .

Thus, under suitable orientation assumptions, Σ admits two smooth nonvanishing future directed null normal vector fields l_+ and l_- (unique up to positive rescaling).



By convention, we refer to l_+ as outward pointing and l_- as inward pointing.

Associated to l_+ and l_- , are the two *null second fundamental forms*, χ_+ and χ_- , respectively, defined as

$$\chi_\pm : T_p\Sigma \times T_p\Sigma \rightarrow \mathbb{R}, \quad \chi_\pm(X, Y) = g(\nabla_X l_\pm, Y). \quad (7.22)$$

The *null expansion scalars* (or *null mean curvatures*) θ_\pm of Σ are obtained by tracing χ_\pm with respect to the induced metric γ on Σ ,

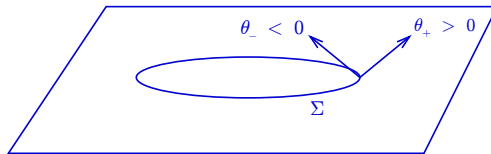
$$\theta_\pm = \text{tr}_\gamma \chi_\pm = \gamma^{AB} \chi_{\pm AB} = \text{div}_\Sigma l_\pm. \quad (7.23)$$

It can be seen that the sign of θ_\pm does not depend on the scaling of l_\pm . Physically, θ_+ (resp., θ_-) measures the divergence of the outgoing (resp., ingoing) light rays emanating orthogonally from Σ .

It is useful to note the connection between the null expansion scalars θ_\pm and the expansion of the generators of a null hypersurface, as discussed in Section 6. Let S_+

be the null hypersurface, defined and smooth near Σ , generated by the null geodesics passing through Σ with initial tangents l_+ . Then θ_+ is the null expansion of S_+ restricted to Σ ; θ_- may be described similarly.

For round spheres in Euclidean slices of Minkowski space, with the obvious choice of inside and outside, one has $\theta_- < 0$ and $\theta_+ > 0$.



In fact, this is the case in general for large “radial” spheres in *asymptotically flat* spacelike hypersurfaces. However, in regions of spacetime where the gravitational field is strong, one may have both $\theta_- < 0$ and $\theta_+ < 0$, in which case Σ is called a *trapped surface*. As we now discuss, under appropriate energy and causality conditions, the occurrence of a trapped surface signals the onset of gravitational collapse [16]. See [4, 13] for results concerning the dynamical formation of trapped surfaces.

The Penrose singularity theorem [16] is the first of the famous singularity theorems of general relativity. The singularity theorems establish, under generic circumstances, the existence in spacetime of incomplete timelike or null geodesics. Such incompleteness indicates that spacetime has come to an end either in the past or future. In specific models past incompleteness is typically associated with a “big bang” beginning of the universe, and future incompleteness is typically associated with a “big crunch” (time dual of the big bang), or, of a more local nature, gravitational collapse to a black hole. The Penrose singularity theorem is associated with the latter.

All the classical singularity theorems require energy conditions. The Penrose singularity theorem requires that $\text{Ric}(X, X) \geq 0$ for all null vectors X . Note that for spacetimes satisfying the Einstein equations, this is just the null energy condition (NEC), cf. Equation (1.13).

In studying an *isolated* gravitating system, such as the collapse of a star and formation of a black hole, it is customary to model this situation by a spacetime which is asymptotically flat (i.e., asymptotically Minkowskian). In this context, the assumption of the Penrose singularity theorem that spacetime admit a noncompact Cauchy surface is natural.

The key concept introduced by Penrose in this singularity theorem is that of the trapped surface, as discussed above. What Penrose proved is that once the gravitational field becomes sufficiently strong that trapped surfaces appear (as they do in the Schwarzschild solution) then the development of singularities is inevitable.

Theorem 7.1. *Let M be a globally hyperbolic spacetime with noncompact Cauchy surfaces satisfying the NEC. If M contains a trapped surface Σ then M is future null geodesically incomplete.*

Suppose that M is future null geodesically complete. We show that the achronal boundary $\partial I^+(\Sigma)$ is compact. Since $\partial I^+(\Sigma)$ is closed, if $\partial I^+(\Sigma)$ is noncompact, there exists a sequence of points $\{q_n\} \subset \partial I^+(\Sigma)$ that *diverges to infinity* in M , i.e., that does not have a convergent subsequence in M . Since, by Proposition 4.3, $J^+(\Sigma)$ is closed, we have,

$$\partial I^+(\Sigma) = \partial J^+(\Sigma) = J^+(\Sigma) \setminus I^+(\Sigma). \quad (7.24)$$

Hence, by Proposition 2.1, there exists a future directed null geodesic $\eta_n; [0, a_n] \rightarrow M$ from some point $p_n \in \Sigma$ to q_n , which is contained in $\partial I^+(\Sigma)$. In particular, η_n must meet Σ orthogonally at p_n (otherwise $q_n \in I^+(\Sigma)$, cf. [15, Lemma 50, p. 298]).

Since Σ is compact there exists a subsequence $\{p_m\}$ of $\{p_n\}$, such that $p_m \rightarrow p \in \Sigma$. It follows that the sequence $\{\eta_m\}$ converges in the sense of geodesics to a future complete null normal geodesic $\eta : [0, \infty) \rightarrow M$, starting at p , which is contained in $\partial I^+(\Sigma)$. Without loss of generality we may assume η is outward pointing, i.e., $\eta'(0) = l_+(p)$. By Equation (7.24), there can be no timelike curve from a point of Σ to a point of η . This implies that no outward pointing null normal geodesic can meet η , for they would have to meet in a corner. A point further out on η would then be timelike related to Σ . On similar grounds, there can be no *null focal point* to Σ along η , i.e., no point on η where nearby outward pointing null normal geodesics cross η “to first order” ([15, Prop. 48, p. 296]). This implies that the exponential map, restricted to the null normal bundle of Σ , is nonsingular along η (see [15], Prop. 30, p. 283 and Cor. 40, p. 290). It follows that for any $a > 0$, the segment $\eta|_{[0, a]}$, is contained in a smooth null hypersurface S , generated by the outward pointing null normal geodesics emanating from a sufficiently small neighborhood of p in Σ . Since Σ is a trapped surface, $\theta^+(p) < 0$. Choose $a > \frac{n-1}{|\theta^+(p)|}$.

Let $s \rightarrow \theta(s)$ be the null mean curvature of S along η . By assumption, $\theta(0) = \theta^+(p) < 0$. As in the proof of Proposition 6.4, the Raychaudhuri equation (6.19) and the NEC imply the differential inequality (6.21), from which it follows that $\theta \rightarrow -\infty$ in an affine parameter time $\leq \frac{n-1}{|\theta^+(p)|} < a$, contradicting the smoothness of S in a neighborhood of $\eta|_{[0, a]}$.

Thus we have shown that if M is future null geodesically complete then $\partial I^+(\Sigma)$ is compact. It now follows from Propositions 3.1 and 4.8 that $\partial I^+(\Sigma)$ is a *compact* Cauchy surface for M , contrary assumption. \square

For certain applications, the following variant of the Penrose singularity theorem is useful.

Theorem 7.2. *Let M be a globally hyperbolic spacetime satisfying the null energy condition, with smooth spacelike Cauchy surface V . Let Σ be a smooth closed (compact without boundary) hypersurface in V which separates V into an “inside” U and an “outside” W , i.e., $V \setminus \Sigma = U \cup W$ where $U, W \subset V$ are connected disjoint sets. Suppose, further, that \bar{W} is non-compact. If Σ is **outer-trapped** ($\theta_+ < 0$) then M is future null geodesically incomplete.*



Proof. Exercise. Hint: Consider the achronal boundary $\partial I^+(\bar{U})$ and argue similarly to the proof of the Penrose singularity theorem.

This version of the Penrose singularity theorem may be used to prove the following beautiful result of Gannon [10] and Lee [14].

Theorem 7.3. *Let M be a globally hyperbolic spacetime which satisfies the null energy condition and which contains a smooth asymptotically flat spacelike Cauchy surface V . If V is not simply connected ($\pi_1(V) \neq 0$) then M is future null geodesically incomplete.*

Thus, as suggested by this theorem, nontrivial topology tends to induce gravitational collapse. In the standard collapse scenario (based on the *weak cosmic censorship conjecture*) the process of gravitational collapse leads to the formation of an event horizon which shields the singularities from view. According to the *principle of topological censorship* the nontrivial topology that induced collapse should end up behind hidden the event horizon, and the region outside the black hole should have simple topology. There are a number of results supporting this view. See [8] for further discussion, relevant references and related results.

Exercise: Let M be a globally hyperbolic spacetime which satisfies the null energy condition and which contains a smooth asymptotically flat spacelike Cauchy surface V . Use Theorem 7.2 to show that if V has more than one asymptotically flat end then M is future null geodesically incomplete. Thus, in Theorem 7.3 one might as well assume that V has only one asymptotically flat end.

References

- [1] J. K. Beem, P. E. Ehrlich, and K. L. Easley, *Global Lorentzian geometry*, second ed., Monographs and Textbooks in Pure and Applied Mathematics, vol. 202, Marcel Dekker Inc., New York, 1996.
- [2] A. N. Bernal and M. Sánchez, *Smoothness of time functions and the metric splitting of globally hyperbolic spacetimes*, *Comm. Math. Phys.* **257** (2005), no. 1, 43–50.
- [3] ———, *Globally hyperbolic spacetimes can be defined as ‘causal’ instead of ‘strongly causal’*, *Classical Quantum Gravity* **24** (2007), no. 3, 745–749.
- [4] D. Christodoulou, *The formation of black holes in general relativity*, *Geometry and analysis. No. 1, Adv. Lect. Math. (ALM)*, vol. 17, Int. Press, Somerville, MA, 2011, pp. 247–283.
- [5] P. T. Chruściel, *Elements of causal theory*, 2011, arXiv:1110.6706v1.
- [6] P. T. Chruściel, E. Delay, G. J. Galloway, and R. Howard, *Regularity of horizons and the area theorem*, *Ann. Henri Poincaré* **2** (2001), no. 1, 109–178.
- [7] P. T. Chruściel, D. E. Grant, and E. Minguzzi, *On differentiability of volume time functions*, 2013, arXiv:1301.2909v1.
- [8] M. Eichmair, G. J. Galloway, and D. Pollack, *Topological censorship from the initial data point of view*, *J. Differential Geom.* **95** (2013), no. 3, 389–405.
- [9] G. J. Galloway and C. Vega, *Achronal limits, lorentzian spheres, and splitting*, 2012, arXiv:1211.2460v2, to appear in *Annales Henri, Poincaré*.
- [10] D. Gannon, *Singularities in nonsimply connected space-times*, *J. Mathematical Phys.* **16** (1975), no. 12, 2364–2367.
- [11] R. Geroch, *Domain of dependence*, *J. Mathematical Phys.* **11** (1970), 437–449.
- [12] S. W. Hawking and G. F. R. Ellis, *The large scale structure of space-time*, Cambridge University Press, London, 1973, Cambridge Monographs on Mathematical Physics, No. 1.
- [13] S. Klainerman, J. Luk, and I. Rodnianski, *A fully anisotropic mechanism for formation of trapped surfaces in vacuum*, 2013, arXiv:1302.5951v1.
- [14] C. W. Lee, *A restriction on the topology of Cauchy surfaces in general relativity*, *Comm. Math. Phys.* **51** (1976), no. 2, 157–162.

- [15] Barrett O'Neill, *Semi-Riemannian geometry*, Pure and Applied Mathematics, vol. 103, Academic Press Inc. [Harcourt Brace Jovanovich Publishers], New York, 1983, With applications to relativity.
- [16] R. Penrose, *Gravitational collapse and space-time singularities*, Phys. Rev. Lett. **14** (1965), 57–59. MR MR0172678 (30 #2897)
- [17] ———, *Techniques of differential topology in relativity*, Society for Industrial and Applied Mathematics, Philadelphia, Pa., 1972, Conference Board of the Mathematical Sciences Regional Conference Series in Applied Mathematics, No. 7. MR MR0469146 (57 #8942)
- [18] Robert M. Wald, *General relativity*, University of Chicago Press, Chicago, IL, 1984.