# THE CRAIG INTERPOLATION LEMMA

## BURTON DREBEN and HILARY PUTNAM

Let $A$ and $B$ be schemata of quantification theory. Assume that $A$ and $B$ are neither valid nor inconsistent and that $A$ implies $B$. *The Craig Interpolation Lemma* asserts that in this case there exists a schema $M$ *containing only predicate letters common to A and B* such that $A$ implies $M$ and $M$ implies $B$. The purpose of the present note is to present a short and elementary proof. Familiarity with Quine's *Methods of Logic* (revised edition, including the appendix) is presumed on the part of the reader.

1. *The Interpolation Lemma for propositional calculus.* It is well known that the analogous "interpolation lemma" for propositional calculus is trivial. The following elegant demonstration is due to Kreisel:

Let $A$ and $B$ be schemata of propositional calculus such that $A$ implies $B$. Let $A = A(p_1, \ldots, p_n)$. Suppose some $p_i$, say $p_1$, does not occur in $B$. Then $A(p_1, p_2, \ldots, p_n) \supset B$ is valid; hence $A(\top, p_2, \ldots, p_n) \supset B$ is valid and $A(\bot, p_2, \ldots, p_n) \supset B$ is valid.

But $A(p_1, p_2, \ldots, p_n) \supset .A(\top, p_2, \ldots, p_n) \lor A(\bot, p_2, \ldots, p_n)$ is valid. Hence $A(p_1, p_2, \ldots, p_n)$ implies $A(\top, p_2, \ldots, p_n) \lor A(\bot, p_2, \ldots, p_n)$ which implies $B$. The "interpolation formula" $A(\top, p_2, \ldots, p_n) \lor A(\bot, p_2, \ldots, p_n)$ may reduce to just $\top$, but in this case $B$, being implied by $\top$, must be a tautology. Also, the "interpolation formula" may reduce to just $\bot$, but in this case $A$, implying $\bot$, must be inconsistent. So, if $A$ and $B$ are neither tautologous nor inconsistent, the "interpolation formula" must reduce to neither $\top$ nor $\bot$, but to a well formed formula of propositional calculus in the usual sense.[1] One of the remaining letters, $p_2, p_3, \ldots, p_n$ may fail to occur in $B$, but repetition of the method must eventually lead to an "interpolation formula" containing only such $p_i$ as occur in $B$; hence, $p_i$ common to both $A$ and $B$.

2. *The Interpolation Lemma for quantification theory.* Assume $A$ implies $B$ and that $A$ and $B$ are neither valid nor inconsistent. Let $A$ be in prenex normal form, and let $\overline{B}$ be a prenex equivalent of $-B$. Since $A$ implies $B$, $A.-B$ is inconsistent, and hence the set consisting of the two schemata $A$ and $\overline{B}$ is inconsistent. But then, by the appendix of *Methods*, there is a deduction, say,

```
*      (1)    A
**     (2)    B̄
**            .
**            .
**            .
**            .
```

Figure 1

with the properties:

(i) Every line except (1) and (2) comes from a preceding line by either UI or EI.

(ii) The conjunction of all the quantifier-free lines in the deduction is truth-functionally inconsistent.

We shall assume, without loss of generality, that Quine's "rule of thumb" (see *Methods*, p. 164) is conformed to in Fig. 1, i.e., each variable "flagging" a given line is alphabetically later than all the free variables in the line to which the given line is subjoined. (This is equivalent to saying that the "flagged" variable is alphabetically later than all the *other* free variables in the line it flags.) And, of course, no variable is "flagged" more than once; otherwise Fig. 1 would not be a deduction.

Let $C_A$ be the set of all and only the quantifier-free lines in Fig. 1 which come ultimately from $A$ (this is relative to some particular analysis of the deduction shown in Fig. 1, of course, but we shall assume a fixed analysis for the purposes of our argument), and let $C_{\bar{B}}$ be the set of all and only the quantifier-free lines which come ultimately from $\bar{B}$. (Every quantifier-free line in Fig. 1 belongs to either $C_A$ or $C_{\bar{B}}$ and not to both, since the rules UI and EI are the only rules used in Fig. 1, and both rules are one-premise rules.) Let the conjunction of the formulas in $C_A$ be $F_A$ and let the conjunction of the formulas in $C_{\bar{B}}$ be $F_{\bar{B}}$. Since $C_A \cup C_{\bar{B}}$ is the set of all quantifier-free lines in Fig. 1, and this is truth functionally inconsistent, the conjunction $F_A \cdot F_{\bar{B}}$ must also be inconsistent, i.e., $F_A$ implies $-F_{\bar{B}}$ truth-functionally.

By the Interpolation Lemma for propositional calculus, there exists a schema $M$ containing only truth-functional constituents common to $F_A$ and $-F_{\bar{B}}$ such that $F_A$ implies $M$ which implies $-F_{\bar{B}}$. The schema $M$ can contain only predicate letters common to $F_A$ and $-F_{\bar{B}}$, and hence only predicate letters common to $A$ and to $B$ (or $\bar{B}$). Moreover, $M$ can contain only individual variables common to $F_A$ and $-F_{\bar{B}}$, and hence common to $C_A$ and $C_{\bar{B}}$.

Let $(Q)$ be the following string of quantifiers: $(Q_1 v_1)(Q_2 v_2) \ldots (Q_n v_n)$, where $v_1, v_2, \ldots, v_n$ are, in alphabetical order, all the variables common to $F_A$ and $-F_{\bar{B}}$, and $Q_i$ is a universal quantifier (an existential quantifier) if $v_i$ was never introduced (was introduced once) by EI in the course of deriving some line belonging to $C_A$. We claim:

$(Q)M$ is the desired interpolation formula, i.e., $A$ implies $(Q)M$ which implies $B$.

*Proof that A implies* $(Q)M$. By just deleting $\overline{B}$ and all steps which come ultimately from $\overline{B}$ in Fig. 1 we obtain a deduction:

$$
\begin{array}{ccc}
* & (1) & A \\
* & (2) & \underline{\phantom{xx}} \\
* & (3) & \underline{\phantom{xx}} \\
* & & \cdot \\
* & & \cdot \\
* & & \cdot \\
* & (k) & \underline{\phantom{xx}}
\end{array}
$$

Figure 2

whose quantifier-free lines are exactly $C_A$. In this deduction the same variables occuring in $C_A$ are flagged as in Fig. 1, and the "rule of thumb" is, of course, still conformed to.

Since $F_A$ implies $M$ truth-functionally, we can obtain $M$ in one TF step.

But then by $n$ uses of UG and EG we can obtain $(Q)M$ as shown in Fig. 3:

$$
\begin{array}{cclll}
* & (1) & A \\
* & (2) & \underline{\phantom{xx}} \\
* & (3) & \underline{\phantom{xx}} \\
* & & \cdot \\
* & & \cdot \\
* & & \cdot \\
* & (k) & \underline{\phantom{xx}} \\
* & (k+1) & \overline{M} & & (1),(2),\ldots,(k),\ \text{TF} \\
* & & \cdot \\
* & & \cdot \\
* & & \cdot \\
* & (k+n+1) & (Q)M
\end{array}
$$

$\left.\phantom{\begin{array}{c}\cdot\\\cdot\\\cdot\end{array}}\right\}$ (UG and EG steps)

Figure 3

That the deduction shown in Fig. 3 conforms to Quine's rules may be checked by observing that whenever a variable $v_i$ is universally generalized in Fig. 3 it was *not* flagged in Fig. 2, by the definition of the string of quantifiers $(Q)$, and so no variable is flagged twice in Fig. 3. Also, the free variables in the line flagged by $v_i$ are just $v_1, v_2, \ldots, v_{i-1}$ because the quantifiers $(Q)$ have to be put on in the order $(Q_n v_n)$, $(Q_{n-1} v_{n-1})$, $(Q_{n-2} v_{n-2}), \ldots$ in Fig. 3, if they are to end up in the order $(Q_1 v_1)(Q_2 v_2)\ldots$ $(Q_n v_n)$ and $v_1, v_2, \ldots, v_{i-1}$ are all alphabetically earlier than $v_i$. So the "rule of thumb" is conformed to. Thus Fig. 3 exhibits a finished deduction of $(Q)M$ from $A$, and hence, by the soundness of Quine's system, $A$ implies $(Q)M$.

*Proof that* $(Q)M$ *implies* $B$. By just deleting $A$ and all steps which come ultimately from $A$ in Fig. 1 we obtain a deduction:

$$
\begin{array}{lll}
* & (1) & \bar{B} \\
* & (2) & \text{---} \\
* & & \cdot \\
* & & \cdot \\
* & & \cdot \\
* & (j) & \text{---}
\end{array}
$$

Figure 4

whose quantifier-free lines are exactly $C_{\bar{B}}$.

If we now adjoin $(Q)M$ as a new premise, we can deduce $M$ by UI and EI in $n$ steps as shown in Fig. 5.

$$
\begin{array}{lll}
* & (1) & \bar{B} \\
* & (2) & \text{---} \\
* & (3) & \text{---} \\
* & & \cdot \\
* & & \cdot \\
* & & \cdot \\
* & (j) & \text{---} \\
** & (j+1) & (Q)M \\
** & & \cdot \\
** & & \cdot \left.\vphantom{\begin{array}{c}.\\.\\.\end{array}}\right\} \quad \text{(UI and EI steps)} \\
** & & \cdot \\
** & (j+n+1) & M
\end{array}
$$

Figure 5

That the deduction shown in Fig. 5 conforms to Quine's rules may be checked by observing that whenever a variable $v_i$ is existentially instantiated in Fig. 5 it was flagged in Fig. 2, by the definition of the string of quantifiers $(Q)$, and hence *not* flagged in Fig. 4 (since Fig. 2 and Fig. 4 can have no line in common). So no variable is flagged twice in Fig. 5. Also, the free variables other than $v_i$ in the line flagged by $v_i$ are just $v_1, v_2, \ldots, v_{i-1}$ and these are all alphabetically earlier than $v_i$. So the "rule of thumb" is conformed to. Thus Fig. 5 exhibits a deduction, though not a finished one, with premises $\bar{B}$ and $(Q)M$.

Since all the lines in $C_{\bar{B}}$ occur in Fig. 5, and $M$ implies $-F_{\bar{B}}$, Fig. 5 contains a truth-functionally inconsistent set of quantifier-free lines, namely $C_{\bar{B}} \cup \{M\}$. Hence the deduction shown in Fig. 5 can be finished by simply adjoining $B$, by a single TF step:

$$
\begin{array}{llll}
** & (j+n+2) & B & (1), (2), \ldots, (j+n+1), \text{TF}
\end{array}
$$

By the soundness of Quine's system, $\bar{B}$ and $(Q)M$ together imply $B$ and hence $(Q)M$ implies $B$, q.e.d.

If we obtain the "interpolation formula" $M$ by Kreisel's method, it may happen that $M$ reduces to the truth-value $\top$, but in this case $B$, being implied by $(Q)\top$, has to be a valid formula; and similarly if $M$ reduces to the truth-value $\bot$, then $A$, implying $(Q)\bot$, has to be a contradictory formula. So if $A$ and $B$ are neither valid nor contradictory $M$ does not reduce to $\top$ or $\bot$, i.e., $M$ is a well formed formula in the usual sense.[1]

We are grateful to W. V. Quine for helpful comments.

## FOOTNOTE

1. *Methods of Logic* (1959) gives a method for eliminating $\top$ and $\bot$ from any formula of the propositional calculus which is neither valid nor contradictory.

*Harvard University*
*Cambridge, Massachusetts*